



# Adaptation to the Range in $K$ -Armed Bandits

Hédi Hadiji, Gilles Stoltz

## ► To cite this version:

| Hédi Hadiji, Gilles Stoltz. Adaptation to the Range in  $K$ -Armed Bandits. 2020. hal-02794382v2

**HAL Id: hal-02794382**

**<https://hal.science/hal-02794382v2>**

Preprint submitted on 10 Nov 2020 (v2), last revised 9 Jun 2022 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Adaptation to the Range in $K$ -Armed Bandits

Hédi Hadiji & Gilles Stoltz

Université Paris-Saclay, CNRS, Laboratoire de mathématiques d’Orsay, 91405, Orsay, France

`hedi.hadiji,gilles.stoltz@math.u-psud.fr`

November 10, 2020

---

## Abstract

We consider stochastic bandit problems with  $K$  arms, each associated with a distribution supported on a given finite range  $[m, M]$ . We do not assume that the range  $[m, M]$  is known and show that there is a cost for learning this range. Indeed, a new trade-off between distribution-dependent and distribution-free regret bounds arises, which, for instance, prevents from simultaneously achieving the typical  $\ln T$  and  $\sqrt{T}$  bounds. For instance, a  $\sqrt{T}$  distribution-free regret bound may only be achieved if the distribution-dependent regret bounds are at least of order  $\sqrt{T}$ . We exhibit a strategy achieving the rates for regret indicated by the new trade-off.

Area of review: Machine learning and data science

OR/MS Subject Classification—Computer science: artificial intelligence; Decision analysis: sequential

Keywords: multiarmed bandits; adversarial learning; cumulative regret; information-theoretic proof techniques

## 1. Introduction

Stochastic multi-armed bandits form a standard setting to deal with sequential decision-making problems like the design of clinical trials (one of the first applications mentioned) online advertisement, online revenue management, queuing and scheduling, etc. (more recent applications that belong to the operations research area).

However virtually all articles on stochastic  $K$ -armed bandits (notable exceptions are discussed below) either assume that distributions of the arms belong to some parametric family (often, one-dimensional exponential families) or are sub-Gaussian with a known parameter  $\sigma^2$ . Among the latter category, the case of the non-parametric family of distributions supported on a known range  $[m, M]$  is of particular interest to us.

We show that the knowledge of the range  $[m, M]$  is a crucial information and that facing bounded bandit problems but ignoring the bounds  $m$  and  $M$  is much harder. We do so by studying what may be achieved and what cannot be achieved anymore when the range  $[m, M]$  is unknown and the strategies need to learn it. We call this problem adaptation to the range, or scale-free regret minimization.

We prove that adaptation to the range is actually possible but that it has a cost: our most striking result (in Section 2.2) is a trade-off between the scale-free distribution-dependent and distribution-free regret bounds that may be achieved. For instance, no strategy adaptive to the range can simultaneously achieve distribution-dependent regret bounds of order  $\ln T$  and distribution-free regret bounds of order  $\sqrt{T}$  (up to polynomial factors), as simple strategies like UCB strategies (by Auer et al., 2002a) do in the case of a known range. Our general trade-off indicates, for instance, that if one wants to keep the

same  $\sqrt{T}$  order of magnitude for the scale-free distribution-free regret bounds, then the best scale-free distribution-dependent rate that may be achieved is  $\sqrt{T}$ .

We also provide (in Section 4) a strategy, based on exponential weights, that adapts to the range and obtains optimal distribution-dependent and distribution-free regret bounds as indicated by the trade-off: these are of respective orders  $T^{1-\alpha}$  and  $T^\alpha$ , where  $\alpha \in [1/2, 1)$  is a parameter of the strategy.

## Literature Review

Optimal scale-free regret minimization under full monitoring is offered by the AdaHedge strategy by De Rooij et al. [2014], which we will use as a building block in Section 4. The main difficulty in adaptation to the range for stochastic bandits is the adaptation to the upper end  $M$  (see Section 5); this is why Honda and Takemura [2015] could provide optimal  $\ln T$  distribution-dependent regret bounds for payoffs lying in ranges of the form  $(-\infty, M]$ , with a known  $M$ . Lattimore [2017] considers models of distributions with a known bound on their kurtosis (a scale-free measure of the skewness of the distributions) and provides a scale-free algorithm based on the median-of-means estimators, with  $\ln T$  distribution-dependent regret bounds. However, bounded bandits can have an arbitrarily high kurtosis, so our settings are not directly comparable (and we think that bounded distributions with an unknown range is a more natural assumption). Cowan and Katehakis [2015] study adaptation to the range but in the restricted case of uniform distributions over unknown intervals; they provide optimal  $\ln T$  distribution-dependent regret bounds for that specific model (the cost for adaptation is mild and lies only in the multiplicative constant before the  $\ln T$ ). See also similar results by Cowan et al. [2018] for Gaussian distributions with unknown means and variances. Additional important references performing adaptation in some sense for (stochastic and adversarial)  $K$ -armed bandits are discussed below.

**Adaptation to the effective range in adversarial bandits.** Gerchinovitz and Lattimore [2016] show that it is impossible to adapt to the so-called effective range in adversarial bandits. A sequence of rewards has effective range smaller than  $b$  if for all rounds  $t$ , rewards  $y_{t,a}$  at this round all lie in an interval of the form  $[m_t, M_t]$  with  $M_t - m_t \leq b$ . The lower bound they exhibit relies on a sequence of changing intervals of fixed size. This problem is thus different from our setting. See also positive results (upper bounds) by Cesa-Bianchi and Shamir [2018] for adaptation to the effective range.

**Adaptation to the variance.** Audibert et al. [2009] consider a variant of UCB called UCB-V, which adapts to the unknown variance. Its analysis assumes that rewards lie in a known range  $[0, M]$ . The results crucially use Bernstein’s inequality (see, for instance, Reminder 3 in Appendix C for a statement of the latter); as Bernstein’s inequality holds for random variables with supports in  $[-\infty, M]$ , the analysis of UCB-V might perhaps be extended to this case as well. Deviation bounds in Bernstein’s inequality contain two terms, a main term scaling with the standard deviation, and a remainder term, scaling with  $M$ . This remainder term, which seems harmless, is a true issue when  $M$  is not known, as indicated by the results of the present article.

**Other criteria.** Wei and Luo [2018], Zimmert and Seldin [2019], Bubeck et al. [2018], and many more, provide strategies for adversarial bandits with rewards in a known range, say  $[0, 1]$ , and adapting to additional regularity in the data, like small variations or stochasticity of the data—but never to the range itself.

## 2. Setting and Main Results

We consider finitely-armed stochastic bandits with bounded and possibly signed rewards. More precisely,  $K \geq 2$  arms are available; we denote by  $[K]$  the set  $\{1, \dots, K\}$  of these arms. With each arm  $a$  is associated a probability distribution  $\nu_a$  lying in some known model  $\mathcal{D}$ ; a model is a set

of probability distributions over  $\mathbb{R}$  with a first moment. The models of interest in this article are discussed below. A bandit problem in  $\mathcal{D}$  is a  $K$ -vector of probability distributions in  $\mathcal{D}$ : we denote it by  $\underline{\nu} = (\nu_a)_{a \in [K]}$ . The player knows  $\mathcal{D}$  but not  $\underline{\nu}$ . As is standard in this setting, we denote by  $\mu_a = \mathbb{E}(\nu_a)$  the mean payoff provided by an arm  $a$ . An optimal arm and the optimal mean payoff are respectively given by  $a^* \in \operatorname{argmax}_{a \in [K]} \mu_a$  and  $\mu^* = \max_{a \in [K]} \mu_a$ . Finally,  $\Delta_a = \mu^* - \mu_a$  denotes the gap of an arm  $a$ .

The online learning game goes as follows: at round  $t \geq 1$ , the player picks an arm  $A_t \in [K]$ , possibly at random according to a probability distribution  $p_t = (p_{t,a})_{a \in [K]}$  based on an auxiliary randomization  $U_{t-1}$ , and then receives and observes a reward  $Z_t$  drawn independently at random according to the distribution  $\nu_{A_t}$ , given  $A_t$ . More formally, a strategy of the player is a sequence of mappings from the observations to the action set,  $(U_0, Z_1, U_1, \dots, Z_{t-1}, U_{t-1}) \mapsto A_t$ , where  $U_0, U_1, \dots$  are i.i.d. random variables independent from all other random variables and distributed according to a uniform distribution over  $[0, 1]$ . At each given time  $T \geq 1$ , we measure the performance of a strategy through its expected regret:

$$R_T(\underline{\nu}) = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T Z_t \right] = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T \mu_{A_t} \right] = \sum_{a=1}^K \Delta_a \mathbb{E}[N_a(T)], \quad (1)$$

where we used the tower rule for the first equality and defined  $N_a(T)$  as the number of times arm  $a$  was pulled between time rounds 1 and  $T$ .

Doob's optional skipping (see Doob, 1953, Chapter III, Theorem 5.2, page 145 for the original reference, see also Chow and Teicher, 1988, Section 5.3 for a more recent reference) indicates that we may assume that i.i.d. sequences of rewards  $(Y_{t,a})_{t \geq 1}$  are drawn beforehand, independently at random, for each arm  $a$  and that the obtained payoff at round  $t \geq 1$  given the choice  $A_t$  equals  $Z_t = Y_{t,A_t}$ . We will use this second formulation in the rest of the paper as it is the closest to the one of oblivious individual sequences described later in Section 4.1.

**Model: bounded signed rewards with unknown range.** For a given range  $[m, M]$ , where  $m < M$  are two real numbers (not necessarily nonnegative), we denote by  $\mathcal{D}_{m,M}$  the set of probability distributions supported on  $[m, M]$ . Then, the model corresponding to distributions with a bounded but unknown range is the union of all such  $\mathcal{D}_{m,M}$ :

$$\mathcal{D}_{-,+} = \bigcup_{m, M \in \mathbb{R}: m < M} \mathcal{D}_{m,M}.$$

## 2.1. Adaptation to the Range: Scale-Free Regret Bounds

Regret scales with the range length  $M - m$ , thus regret bounds involve a multiplicative factor  $M - m$ . We therefore consider such bounds divided by the scale factor  $M - m$  and call them scale-free regret bounds. We denote by  $\mathbb{N}$  the set of natural integers; (rates on) regret bounds will be given by functions  $\Phi : \mathbb{N} \rightarrow [0, +\infty)$ . We define adaptation to the unknown range in Definitions 1 and 2 below.

**Definition 1** (Scale-free distribution-free regret bounds). *A strategy for stochastic bandits is adaptive to the unknown range of payoffs with a scale-free distribution-free regret bound  $\Phi_{\text{free}} : \mathbb{N} \rightarrow [0, +\infty)$  if for all real numbers  $m < M$ , the strategy ensures, without the knowledge of  $m$  and  $M$ :*

$$\forall \underline{\nu} \text{ in } \mathcal{D}_{m,M}, \quad \forall T \geq 1, \quad R_T(\underline{\nu}) \leq (M - m) \Phi_{\text{free}}(T).$$

We show in Section 4 that adaptation to the unknown range may indeed be performed in the sense of Definition 1, with a scale-free distribution-free regret bound of order  $\sqrt{KT \ln K}$ . The latter is optimal up to maybe a factor of  $\sqrt{\ln K}$ : Auer et al. [2002b] provided a lower bound  $(1/20) \min\{\sqrt{KT}, T\}$  on the regret of any strategy against individual sequences in  $[0, 1]^K$ , thus for bandit problems in  $\mathcal{D}_{0,1}$ , thus for scale-free distribution-free regret bounds.

**Definition 2** (Distribution-dependent rates for adaptation). *A strategy for stochastic bandits is adaptive to the unknown range of payoffs with a distribution-dependent rate  $\Phi_{\text{dep}} : \mathbb{N} \rightarrow [0, +\infty)$  if for all real numbers  $m < M$ , the strategy ensures, without the knowledge of  $m$  and  $M$ :*

$$\forall \underline{\nu} \text{ in } \mathcal{D}_{m,M}, \quad \limsup_{T \rightarrow +\infty} \frac{R_T(\underline{\nu})}{\Phi_{\text{dep}}(T)} < +\infty.$$

*Put differently, the strategy ensures that  $\limsup R_T(\underline{\nu})/\Phi_{\text{dep}}(T) < +\infty$  for all  $\underline{\nu} \in \mathcal{D}_{-,+}$ .*

Definition 2 does not add much to the classical notion of distribution-dependent rates on regret bounds, as the scale factor  $M - m$  does not appear in the definition; it merely ensures that the strategy is not informed of the range. This lack of information prevents from achieving the typical  $\ln T$  order of magnitude on the regret: all uniformly fast convergent strategies on  $\mathcal{D}_{-,+}$  are such that, for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}_{-,+}$  with at least one suboptimal arm,

$$\liminf_{T \rightarrow +\infty} \frac{R_T(\underline{\nu})}{\ln T} = +\infty. \quad (2)$$

This follows from showing that  $\mathcal{K}_{\text{inf}}(\nu_a, \mu^*, \mathcal{D}_{-,+}) = 0$ , where  $\mathcal{K}_{\text{inf}}$  is some infimum of Kullback-Leibler divergences. (A strategy is said to be uniformly fast convergent on a model  $\mathcal{D}$  if for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}$  and all  $\alpha \in (0, 1]$ , its regret satisfies  $R_T(\underline{\nu})/T^\alpha \rightarrow 0$ ; this is a minimal requirement when studying lower bounds to rule out poor strategies.) However, any rate  $\Phi_{\text{dep}}(T) \gg \ln T$  may be achieved thanks to a simple upper-confidence bound [UCB] strategy. Further details, including proofs of the two claims above, may be found in Appendix A of the supplementary material.

## 2.2. Simultaneous Regret Bounds

When the range  $[m, M]$  of the payoffs is known, it is possible to simultaneously achieve optimal distribution-free bounds (of order  $\sqrt{KT}$ ) and optimal distribution-dependent bounds (of order  $\ln T$  with the optimal constant given some  $\mathcal{K}_{\text{inf}}$ ); see the KL-UCB-switch strategy by Garivier et al. [2019a]. Put differently, when the range of payoffs is known, one can achieve optimal (asymptotic) distribution-dependent regret bounds while not sacrificing finite-time guarantees. Simpler strategies like UCB strategies (see Auer et al., 2002a) also simultaneously achieve regret bounds of similar  $\ln T$  and  $\sqrt{T \ln T}$  orders of magnitude but with suboptimal constants.

Our first main result indicates that getting simultaneously these  $\ln T$  and  $\sqrt{T}$  rates is not possible anymore when the range of payoffs is unknown.

**Theorem 1.** *Any strategy with a scale-free distribution-free regret bound satisfying  $\Phi_{\text{free}}(T) = o(T)$  may only achieve distribution-dependent rates  $\Phi_{\text{dep}}$  for adaptation satisfying  $\Phi_{\text{dep}}(T) \geq T/\Phi_{\text{free}}(T)$ .*

*More precisely, the regret of such a strategy is lower bounded as: for all  $\underline{\nu}$  in  $\mathcal{D}_{-,+}$ ,*

$$\liminf_{T \rightarrow \infty} \frac{R_T(\underline{\nu})}{T/\Phi_{\text{free}}(T)} \geq \frac{1}{16} \sum_{a=1}^K \Delta_a. \quad (3)$$

The orders of magnitude of the scale-free distribution-free regret bounds  $\Phi_{\text{free}}(T)$  range between the optimal  $\sqrt{T}$  and the trivial  $T$  rates. The distribution-dependent rates  $\Phi_{\text{dep}}$  for adaptation to the range are therefore at best  $\sqrt{T}$  for strategies enjoying scale-free distribution-free regret bounds;  $\ln T$  rates are excluded. More generally, Theorem 1 shows that there is a trade-off: to force faster distribution-dependent rates for adaptation, one must suffer worsened scale-free distribution-free regret bounds.

The proof of Theorem 1 is provided in Section 3. It actually provides a finite-time (but messy) lower bound on  $R_T(\underline{\nu})/(T/\Phi_{\text{free}}(T))$ .

Our second main result consists of showing that the trade-off imposed by Theorem 1 may indeed be achieved. Section 4 will introduce a strategy, called AHB (AdaHedge for  $K$ -armed Bandits with

extra-exploration, see Algorithm 1) and relying on a parameter  $\alpha \in [1/2, 1)$ . Theorems 2 and 3 show that AHB adapts to the unknown range, satisfies a scale-free distribution-free regret bound

$$\Phi_{\text{free}}^{\text{AHB}}(T) = \left(3 + \frac{5}{\sqrt{1-\alpha}}\right) \sqrt{K \ln K} T^\alpha + 10K \ln K = \mathcal{O}(T^\alpha)$$

and achieves a distribution-dependent rate for adaptation  $\Phi_{\text{dep}}^{\text{AHB}}(T) = T/\Phi_{\text{free}}^{\text{AHB}}(T) = \mathcal{O}(T^{1-\alpha})$ . Even better, we show that for all  $\underline{\nu}$  in  $\mathcal{D}_{-,+}$ ,

$$\limsup_{T \rightarrow \infty} \frac{R_T(\underline{\nu})}{T/\Phi_{\text{free}}^{\text{AHB}}(T)} \leq \frac{12 \ln K}{1-\alpha} \sum_{a=1}^K \Delta_a. \quad (4)$$

The distribution-dependent constants in the right-hand sides of (3) and (4) differ only by a numerical factor of  $16 \times 12/(1-\alpha)$  and a  $\ln K$  factor.

### 2.3. Linear Bandits

We consider the simplest setting of stochastic bandits in this article: with finitely many arms. However, the techniques developed for adaptation to the range in this setting may be generalized to deal, e.g., with (oblivious) adversarial linear bandits; see details in an online appendix [arXiv:2006.03378], Section D.

## 3. Proof of Theorem 1

We follow a proof technique introduced by Lai and Robbins [1985] and Burnetas and Katehakis [1996] and recently revisited by Garivier et al. [2019b]. We fix some bandit problem  $\underline{\nu}$  in  $\mathcal{D}_{-,+}$  and construct an alternative bandit problem  $\underline{\nu}'$  in  $\mathcal{D}_{-,+}$  by modifying the distribution of a single suboptimal arm  $a$  to make it optimal (which is always possible, as there is no bound on the upper end on the ranges of the payoffs in the model). We apply a fundamental inequality that links the expectations of the numbers of times  $N_a(T)$  that  $a$  is pulled under  $\underline{\nu}$  and  $\underline{\nu}'$ . We then substitute inequalities stemming from the definition of distribution-free scale-free regret bounds  $\Phi_{\text{free}}$ , and the result follows by rearranging all inequalities.

**Step 1: Alternative bandit problem.** The lower bound is trivial (it equals 0) when all arms of  $\underline{\nu}$  are optimal. We therefore assume that at least one arm is suboptimal and fix such an arm  $a$ . For some  $\varepsilon \in [0, 1]$  to be defined later by the analysis, we introduce the alternative problem  $\underline{\nu}' = (\nu'_k)_{k \in [K]}$  with  $\nu'_k = \nu_k$  for  $j \neq a$  and  $\nu'_a = (1-\varepsilon)\nu_a + \varepsilon\delta_{\mu_a+2\Delta_a/\varepsilon}$ . This distribution  $\nu'_a$  has a bounded range, so that  $\underline{\nu}'$  lies indeed in  $\mathcal{D}_{-,+}$ . The expectation of  $\nu'_a$  equals  $\mu'_a = \mu_a + 2\Delta_a = \mu^* + \Delta_a > \mu^*$ . Thus,  $a$  is the only optimal arm in  $\underline{\nu}'$ . Finally, for  $\varepsilon$  small enough,  $\mu_a + 2\Delta_a/\varepsilon$  lies outside of the bounded support of  $\nu_a$ . In that case, the density of  $\nu_a$  with respect to  $\nu'_a$  is given by  $1/(1-\varepsilon)$  on the support of  $\nu_a$  (and 0 elsewhere), so that  $\text{KL}(\nu_a, \nu'_a) = \ln(1/(1-\varepsilon))$ .

**Step 2: Application of a fundamental inequality.** We denote by  $\text{kl}(p, q)$  the Kullback-Leibler divergence between Bernoulli distributions with parameters  $p$  and  $q$ . We also index expectations in the rest of the proof by the bandit problem they are relative to: for instance,  $\mathbb{E}_{\underline{\nu}}$  denotes the expectation of a random variable when the ambient randomness is given by the bandit problem  $\underline{\nu}$ . The fundamental inequality for lower bounds on the regret of stochastic bandits (Garivier et al., 2019b, Section 2, Equation 6), which is based on the chain rule for Kullback-Leibler divergence and on a data-processing inequality for expectations of  $[0, 1]$ -valued random variables, reads:

$$\text{kl}\left(\frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{T}, \frac{\mathbb{E}_{\underline{\nu}'}[N_a(T)]}{T}\right) \leq \mathbb{E}_{\underline{\nu}}[N_a(T)] \text{KL}(\nu_a, \nu'_a) = \mathbb{E}_{\underline{\nu}}[N_a(T)] \ln(1/(1-\varepsilon)).$$

Now, since  $u \in (-\infty, 1) \mapsto -u^{-1} \ln(1-u)$  is increasing, we have  $\ln(1/(1-\varepsilon)) \leq (2 \ln 2)\varepsilon$  for  $\varepsilon \leq 1/2$ . For all  $(p, q) \in [0, 1]^2$  and with the usual measure-theoretic conventions,

$$\text{kl}(p, q) = \underbrace{p \ln p + q \ln q}_{\geq -\ln 2} + \underbrace{p \ln \frac{1}{q}}_{\geq 0} + (1-p) \ln \frac{1}{1-q} \geq (1-p) \ln \frac{1}{1-q} - \ln 2,$$

so that, putting all inequalities together, we have proved

$$\left(1 - \frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{T}\right) \ln\left(\frac{1}{1 - \mathbb{E}_{\underline{\nu}'}[N_a(T)]/T}\right) - \ln 2 \leq (2 \ln 2) \varepsilon \mathbb{E}_{\underline{\nu}}[N_a(T)]. \quad (5)$$

So far, we only imposed the constraint  $\varepsilon \in [0, 1/2]$ .

**Step 3: Inequalities stemming from the definition of scale-free distribution-free regret bounds.** We denote by  $[m, M]$  a range containing the supports of all distributions of  $\underline{\nu}$ . By definition of  $\Phi_{\text{free}}$ , given that  $a$  is a suboptimal arm (i.e.,  $\Delta_a > 0$ ):

$$\Delta_a \mathbb{E}_{\underline{\nu}}[N_a(T)] \leq R_T(\underline{\nu}) \leq (M - m) \Phi_{\text{free}}(T).$$

Because of  $\nu'_a$ , the distributions of  $\underline{\nu}'$  have supports within the range  $[m, M_\varepsilon]$ , where we denoted  $M_\varepsilon = \max\{M, \mu_a + 2\Delta_a/\varepsilon\}$ . For  $\underline{\nu}'$ , by definition of  $\Phi_{\text{free}}$ , and given that all gaps  $\Delta'_k$  are larger than the gap  $\Delta'_a = \mu'_a - \mu^* = \Delta_a$  between the unique optimal  $a$  and the second best arms (which were the optimal arms of  $\underline{\nu}$ ),

$$\Delta_a(T - \mathbb{E}_{\underline{\nu}'}[N_a(T)]) = \Delta'_a(T - \mathbb{E}_{\underline{\nu}'}[N_a(T)]) \leq \sum_{j \neq a} \Delta'_j \mathbb{E}_{\underline{\nu}'}[N_j(T)] = R_T(\underline{\nu}') \leq (M_\varepsilon - m) \Phi_{\text{free}}(T).$$

By rearranging the two inequalities above, we get

$$1 - \frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{T} \geq 1 - \frac{(M - m) \Phi_{\text{free}}(T)}{T \Delta_a} \quad \text{and} \quad 1 - \frac{\mathbb{E}_{\underline{\nu}'}[N_a(T)]}{T} \leq \frac{(M_\varepsilon - m) \Phi_{\text{free}}(T)}{T \Delta_a},$$

thus, after substitution into (5),

$$\left(1 - \frac{(M - m) \Phi_{\text{free}}(T)}{T \Delta_a}\right) \ln\left(\frac{T \Delta_a}{(M_\varepsilon - m) \Phi_{\text{free}}(T)}\right) - \ln 2 \leq (2 \ln 2) \varepsilon \mathbb{E}_{\underline{\nu}}[N_a(T)]. \quad (6)$$

**Step 4: Final calculations.** We take  $\varepsilon = \varepsilon_T = \alpha^{-1} \Phi_{\text{free}}(T)/T$  for some constant  $\alpha > 0$ ; we will pick  $\alpha = 1/8$ . By the assumption  $\Phi_{\text{free}}(T) = o(T)$ , we have  $\varepsilon_T \leq 1/2$  as needed for  $T$  large enough, as well as  $M_{\varepsilon_T} = \mu_a + 2\Delta_a/\varepsilon_T = \mu_a + 2\alpha\Delta_a T/\Phi_{\text{free}}(T)$ . Substituting these values into (6), a finite-time lower bound on the quantity of interest is finally given by

$$\frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{T/\Phi_{\text{free}}(T)} \geq \frac{\alpha}{2 \ln 2} \left( -\ln 2 + \underbrace{\left(1 - \frac{(M - m) \Phi_{\text{free}}(T)}{T \Delta_a}\right)}_{\rightarrow 0} \ln \underbrace{\left(\frac{T \Delta_a}{2\alpha\Delta_a T + (\mu_a - m)\Phi_{\text{free}}(T)}\right)}_{\rightarrow 1/(2\alpha)} \right).$$

It entails the asymptotic lower bound

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{T/\Phi_{\text{free}}(T)} \geq \frac{\alpha}{2 \ln 2} (\ln(1/\alpha) - 2 \ln 2) = \frac{1}{16}$$

for the choice  $\alpha = 1/8$ . The claimed result follows by adding these lower bounds for each suboptimal arm  $a$ , with a factor  $\Delta_a$ , following the rewriting (1) of the regret.

**Remark 1.** *The proof above only exploits the fact that the upper end  $M$  of the range is unknown: the alternative problems lie in  $\mathcal{D}_{m,M'}$  for some  $M'$  that can be arbitrarily large. Yet, by definition of adaptation to the range, the strategy needs to guarantee  $(M' - m) \Phi_{\text{free}}(T)$  distribution-free regret bounds in that case.*

We may note that therefore, Theorem 1 also holds for the model of bounded distributions with a known lower end  $m \in \mathbb{R}$  for the range:

$$\mathcal{D}_{m,+} = \bigcup_{\substack{M \in \mathbb{R}: \\ M > m}} \mathcal{D}_{m,M}. \quad (7)$$

Definitions 1 and 2 handle the case of  $\mathcal{D}_{-,+}$  but can be adapted in an obvious way to  $\mathcal{D}_{m,+}$  by fixing  $m$ , by having the strategy know  $m$ , and require the bounds to hold for all  $M \in [m, +\infty)$  and all bandit problems in  $\mathcal{D}_{m,M}$ , thus leading to the concept of adaptation to the upper end of the range.

## 4. Adaptation to Range Based on AdaHedge: The AHB Strategy

When the range of payoffs is known, Auer et al. [2002b] use exponential weights (Hedge) on estimated payoffs and with extra-exploration (mixing with the uniform distribution) to achieve a distribution-free regret bound of order  $\sqrt{KT \ln K}$ . Actually, it is folklore knowledge that the extra-exploration used in this case is unnecessary. To deal with the case of unknown payoff range, we consider a self-tuned version of Hedge called AdaHedge (De Rooij et al., 2014, see also earlier work by Cesa-Bianchi et al., 2007) and do add extra-exploration. Just as Auer et al. [2002b], we will actually obtain regret guarantees for oblivious adversarial bandits, not only distribution-free regret bounds for stochastic bandits. We therefore introduce now the setting of oblivious adversarial bandits and define adaptation to the range in that case.

### 4.1. Oblivious Adversarial Bandits

In the setting of (fully) oblivious adversarial bandits (see Cesa-Bianchi and Lugosi, 2006, Audibert and Bubeck, 2009), a range  $[m, M]$  is set by the environment, where  $m, M$  are real numbers (not necessarily nonnegative), and the environment picks beforehand a sequence  $y_1, y_2, \dots$  of reward vectors in  $[m, M]^K$ . We denote by  $y_t = (y_{t,a})_{a \in [K]}$  the components of these vectors. The online learning game starts only then: at each round  $t \geq 1$ , the player picks an arm  $A_t \in [K]$ , possibly at random according to a probability distribution  $p_t = (p_{t,a})_{a \in [K]}$  based on an auxiliary randomization  $U_{t-1}$ , and then receives and observes  $y_{t,A_t}$ . More formally, a strategy of the player is a sequence of mappings from the observations to the action set,  $(U_0, y_{1,A_1}, U_1, \dots, y_{t-1,A_{t-1}}, U_{t-1}) \mapsto A_t$ , where  $U_0, U_1, \dots$  are i.i.d. random variables independent from all other random variables and distributed according to a uniform distribution over  $[0, 1]$ . At each given time  $T \geq 1$ , denoting by  $y_{1:T} = (y_1, \dots, y_T)$  the reward vectors, we measure the performance of a strategy through its expected regret:

$$R_T(y_{1:T}) = \max_{a \in [K]} \sum_{t=1}^T y_{t,a} - \mathbb{E} \left[ \sum_{t=1}^T y_{t,A_t} \right], \quad (8)$$

where all randomness lies in the choice of the arms  $A_t$  only (i.e., where the expectation is over the choice of the arms  $A_t$  only) as rewards are fixed beforehand.

The counterpart of Definition 1 in this setting is stated next.

**Definition 3** (Scale-free adversarial regret bounds). *A strategy for oblivious adversarial bandits is adaptive to the unknown range of payoffs with a scale-free adversarial regret bound  $\Phi_{\text{adv}} : \mathbb{N} \rightarrow [0, +\infty)$  if for all real numbers  $m < M$ , the strategy ensures, without the knowledge of  $m$  and  $M$ :*

$$\forall y_1, y_2, \dots \text{ in } [m, M]^K, \quad \forall T \geq 1, \quad R_T(y_{1:T}) \leq (M - m) \Phi_{\text{adv}}(T).$$



**Conversion of upper/lower bounds from one setting to the other.** By the tower rule for the right-most equality, we note that for all  $m < M$  and for all  $\underline{\nu}$  in  $\mathcal{D}_{m,M}$ ,

$$R_T(\underline{\nu}) = \max_{a \in [K]} \mathbb{E} \left[ \sum_{t=1}^T Y_{t,a} \right] - \mathbb{E} \left[ \sum_{t=1}^T Y_{t,A_t} \right] \leq \mathbb{E} \left[ \max_{a \in [K]} \sum_{t=1}^T Y_{t,a} - \sum_{t=1}^T Y_{t,A_t} \right] = \mathbb{E} [R_T(Y_{1:T})] \leq \sup_{y_{1:T} \text{ in } [m,M]^K} R_T(y_{1:T}).$$

In particular, lower bounds on the regret for stochastic bandits are also lower bounds on the regret for oblivious adversarial bandits, and strategies designed for oblivious adversarial bandits obtain the same distribution-free regret bounds for stochastic bandits when the individual payoffs  $y_{t,A_t}$  in their definition are replaced with the stochastic payoffs  $Y_{t,A_t}$ .

## 4.2. The AHB Strategy

We state our main strategy, AHB (AdaHedge for  $K$ -armed Bandits, with extra-exploration), in the setting of oblivious adversarial bandits, see Algorithm 1. In a setting of stochastic bandits, it suffices to replace therein  $y_{t,A_t}$  with  $Y_{t,A_t}$ . The AHB strategy relies on a payoff estimation scheme, which we discuss now.

In Algorithm 1, some initial exploration lasting  $K$  rounds is used to get a rough idea of the location of the payoffs and to center the estimates used at an appropriate location. Following Auer et al. [2002b], we consider, for all rounds  $t \geq K + 1$  and arms  $a \in [K]$ ,

$$\hat{y}_{t,a} = \frac{y_{t,A_t} - C}{p_{t,a}} \mathbb{1}_{\{A_t=a\}} + C \quad \text{where} \quad C \stackrel{\text{def}}{=} \frac{1}{K} \sum_{s=1}^K y_{s,s}. \quad (9)$$

Note that all  $p_{t,a} > 0$  for Algorithm 1 due to the use of exponential weights. As proved by Auer et al. [2002b], the estimates  $\hat{y}_{t,a}$  are (conditionally) unbiased. Indeed, the distributions  $q_t$  and  $p_t$  (as well as the constant  $C$ ) are measurable functions of the information  $H_{t-1} = (U_0, y_{1,A_1}, U_1, \dots, U_{t-2}, y_{t-1,A_{t-1}})$  available at the beginning of round  $t \geq K + 1$ , and the arm  $A_t$  is drawn independently at random according to  $p_t$  based on an auxiliary randomization denoted by  $U_{t-1}$ . Therefore, given that the payoffs are oblivious, the conditional expectation of  $\hat{y}_{t,a}$  with respect to  $H_{t-1}$  amounts to integrating over the randomness given by the random draw  $A_t \sim p_t$ : for  $t \geq K + 1$ ,

$$\mathbb{E}[\hat{y}_{t,a} \mid H_{t-1}] = \frac{y_{t,a} - C}{p_{t,a}} \mathbb{P}(A_t = a \mid H_{t-1}) + C = \frac{y_{t,a} - C}{p_{t,a}} p_{t,a} + C = y_{t,a}. \quad (10)$$

These estimators are bounded: assuming that all  $y_{t,a}$ , thus also  $C$ , belong to the range  $[m, M]$ , and given that the distributions  $p_t$  were obtained by a mixing with the uniform distribution, with weight  $\gamma_t$ , we have  $p_{t,a} \geq \gamma_t/K$ , and therefore,

$$\forall t \geq K + 1, \quad \forall a \in [K], \quad |\hat{y}_{t,a} - C| \leq \frac{|y_{t,a} - C|}{p_{t,a}} \leq \frac{M - m}{\gamma_t/K}. \quad (11)$$

**Remark 2.** Algorithm 1 is invariant by affine changes (translations and/or multiplications by positive factors) of the payoffs, as AdaHedge (see De Rooij et al., 2014, Theorem 16) and the payoff estimation scheme (9) are so. This is key for adaptation to the range.

## 4.3. Regret Analysis, Part 1: Scale-Free Adversarial Regret Bound

**Theorem 2.** AdaHedge for  $K$ -armed bandits (Algorithm 1) with a non-increasing extra-exploration sequence  $(\gamma_t)_{t \geq 1}$  smaller than  $1/2$  and the estimation scheme given by (9) ensures that for all bounded ranges  $[m, M]$ , for all oblivious individual sequences  $y_1, y_2, \dots$  in  $[m, M]^K$ , for all  $T \geq 1$ ,

$$R_T(y_{1:T}) \leq 3(M - m) \sqrt{KT \ln K} + 5(M - m) \frac{K \ln K}{\gamma_T} + (M - m) \sum_{t=K+1}^T \gamma_t.$$

---

**Algorithm 1** AHB: AdaHedge for  $K$ -armed Bandits, with extra-exploration
 

---

- 1: **Input:** a sequence  $(\gamma_t)_{t \geq 1}$  in  $[0, 1]$  of extra-exploration rates; a payoff estimation scheme
- 2: **for** rounds  $t = 1, \dots, K$  **do**
- 3:   Draw arm  $A_t = t$
- 4:   Get and observe the payoff  $y_{t,t}$
- 5: **end for**
- 6: **AdaHedge initialization:**  $\eta_{K+1} = +\infty$  and  $q_{K+1} = (1/K, \dots, 1/K) \stackrel{\text{def}}{=} \mathbf{1}/K$
- 7: **for** rounds  $t = K + 1, \dots$  **do**
- 8:   Define  $p_t$  by mixing  $q_t$  with the uniform distribution according to  $p_t = (1 - \gamma_t)q_t + \gamma_t \mathbf{1}/K$
- 9:   Draw an arm  $A_t \sim p_t$  (independently at random according to the distribution  $p_t$ )
- 10:   Get and observe the payoff  $y_{t,A_t}$
- 11:   Compute estimates  $\hat{y}_{t,a}$  of all payoffs with the payoff estimation scheme considered
- 12:   Compute the mixability gap  $\delta_t \geq 0$  based on the distribution  $q_t$  and on these estimates:

$$\delta_t = - \sum_{a=1}^K q_{t,a} \hat{y}_{t,a} + \frac{1}{\eta_t} \ln \left( \sum_{a=1}^K q_{t,a} e^{\eta_t \hat{y}_{t,a}} \right), \quad \text{with} \quad \underbrace{\delta_t = - \sum_{a=1}^K q_{t,a} \hat{y}_{t,a} + \max_{a \in [K]} \hat{y}_{t,a}}_{\text{when } \eta_t = +\infty}$$

- 13:   Compute the learning rate  $\eta_{t+1} = \left( \sum_{s=K+1}^t \delta_s \right)^{-1} \ln K$
- 14:   Define  $q_{t+1}$  component-wise as

$$q_{t+1,a} = \exp \left( \eta_{t+1} \sum_{s=K+1}^t \hat{y}_{a,s} \right) / \sum_{k=1}^K \exp \left( \eta_{t+1} \sum_{s=K+1}^t \hat{y}_{k,s} \right)$$

- 15: **end for**
-

In particular, given a parameter  $\alpha \in (0, 1)$ , the extra-exploration  $\gamma_t = \min\{1/2, \sqrt{5(1-\alpha)K \ln K}/t^\alpha\}$  leads to the scale-free adversarial regret bound

$$\Phi_{\text{adv}}(T) = \left(3 + \frac{5}{\sqrt{1-\alpha}}\right)(M-m)\sqrt{K \ln K} T^{\max\{\alpha, 1-\alpha\}} + 10(M-m)K \ln K. \quad (12)$$

For  $\alpha = 1/2$ , the bound reads  $\Phi_{\text{adv}}(T) = 7(M-m)\sqrt{TK \ln K} + 10(M-m)K \ln K$ .

This value  $\alpha = 1/2$  is the best one to consider if one is only interested in a distribution-free bound (i.e., one is not interested in the distribution-dependent rates for the regret). The proof of Theorem 2 is detailed in Appendix B but we sketch its proof here.

*Proof sketch.* A direct application of the AdaHedge regret bound (Lemma 3 and Theorem 6 of De Rooij et al., 2014), bounding the variance terms of the form  $\mathbb{E}[(X - \mathbb{E}[X])^2]$  by  $\mathbb{E}[(X - C)^2]$ , ensures that

$$\max_{k \in [K]} \sum_{t=K+1}^T \hat{y}_{t,k} - \sum_{\substack{t \geq K+1 \\ a \in [K]}} q_{t,a} \hat{y}_{t,a} \leq 2 \sqrt{\sum_{\substack{t \geq K+1 \\ a \in [K]}} q_{t,a} (\hat{y}_{t,a} - C)^2 \ln K} + \frac{M-m}{\gamma_T/K} \left(2 + \frac{4}{3} \ln K\right).$$

We take expectations, use the definition of the  $p_t$  in terms of the  $q_t$  in the left-hand side, and apply Jensen's inequality in the right-hand side to get

$$\begin{aligned} \mathbb{E} \left[ \max_{k \in [K]} \sum_{t=K+1}^T \hat{y}_{t,k} - \sum_{t=K+1}^T \sum_{a=1}^K \overbrace{p_{t,a} \hat{y}_{t,a}}^{=y_{t,A_t}} + \sum_{t=K+1}^T \gamma_t \sum_{a=1}^K \overbrace{(1/K - q_{t,a}) \hat{y}_{t,a}}^{\mathbb{E}[\dots] \in [m-M, M-m]} \right] \\ \leq 2 \sqrt{\sum_{t=K+1}^T \sum_{a=1}^K \mathbb{E}[q_{t,a} (\hat{y}_{t,a} - C)^2] \ln K} + \frac{M-m}{\gamma_T/K} \left(2 + \frac{4}{3} \ln K\right). \end{aligned}$$

Since  $p_{t,a} \geq (1-\gamma_t)q_{t,a}$  with  $\gamma_t \leq 1/2$  by assumption on the extra-exploration rate, we have the bound  $q_{t,a} \leq 2p_{t,a}$ . Together with standard calculations similar to (10), we have

$$\mathbb{E}[q_{t,a} (\hat{y}_{t,a} - C)^2] \leq 2 \mathbb{E}[p_{t,a} (\hat{y}_{t,a} - C)^2 \mid H_{t-1}] = 2 \mathbb{E}\left[\frac{(y_{t,A_t} - C)^2}{p_{t,a}} \mathbb{1}_{\{A_t=a\}}\right] = 2 \underbrace{(y_{t,a} - C)^2}_{\leq (M-m)^2}.$$

The proof of the first regret bound of the theorem is concluded by collecting all bounds and by taking care of the first  $K$  rounds. The second regret bound then follows from straightforward calculations.  $\square$

#### 4.4. Regret Analysis, Part 2: Distribution-Dependent Rates for Adaptation

Given the conversion explained in Section 4.1, Algorithm 1 tuned as in Corollary 2 for  $\alpha \in [1/2, 1)$  also enjoys the scale-free distribution-free regret bound  $\Phi_{\text{free}}^{\text{AHB}}(T) = \Phi_{\text{adv}}^{\text{AHB}}(T)$  of order  $T^\alpha$ . The theorem below indicates that AHB is adaptive to the unknown range with a distribution-dependent regret rate  $T/\Phi_{\text{free}}^{\text{AHB}}(T)$  of order  $T^{1-\alpha}$  that is optimal given the lower bound indicated by Theorem 1.

**Theorem 3.** Consider Algorithm 1 tuned with some  $\alpha \in [1/2, 1)$  as in the second part of Theorem 2. For all distributions  $\nu_1, \dots, \nu_K$  in  $\mathcal{D}_{-,+}$ ,

$$\limsup_{T \rightarrow \infty} \frac{R_T(\underline{\nu})}{T/\Phi_{\text{free}}^{\text{AHB}}(T)} \leq \frac{12 \ln K}{1-\alpha} \sum_{a=1}^K \Delta_a. \quad (13)$$

The proof is provided in Appendix C. It follows quite closely that of Theorem 3 in Seldin and Lugosi [2017], where the authors study a variant of the Exp3 algorithm of Auer et al. [2002b] for stochastic rewards. It consists, in our setting, in showing that the number of times the algorithm chooses suboptimal arms is almost only determined by the extra-exploration. Our proof is simpler as we aim for cruder bounds. The main technical difference and issue to solve lies in controlling the learning rates  $\eta_t$ , which heavily depend on data in our case.

## 5. The Case of Bandits with Gains or Losses

It is folklore knowledge that there is a difference in nature between dealing with nonnegative payoffs (gains) or dealing with nonpositive payoffs (losses) for regret minimization under bandit monitoring; see Cesa-Bianchi and Lugosi [2006, Remark 6.5, page 164] for an early reference and Kwon and Perchet [2016] for a more complete literature review. Actually, 0 plays no special role, the issue is rather whether one end of the payoff range is known.

**Known lower end  $m$  on the payoff range.** In that case we deal (up to a translation) with gains. This knowledge does not provide any advantage. Indeed, the impossibility results of Section 2 still hold, namely, no  $\ln T$  rate may be achieved for scale-free distribution-dependent regret bounds (see Theorem 4 in Appendix A) and the same trade-off exists between scale-free distribution-free and distribution-dependent regret bounds (Theorem 1 still holds, see Remark 1).

**Known upper end  $M$  on the payoff range.** What follows is stated with greater details and proved in an online appendix [arXiv:2006.03378], Section E. In that case we deal (up to a translation) with losses, also known as semi-bounded rewards. The DMED strategy of Honda and Takemura [2015] does achieve the optimal asymptotic distribution-dependent regret bound, of order  $\ln T$ . We also show that the INF strategy of Audibert and Bubeck [2009] may be extended, with a little but not too much work, to provide a scale-free distribution-free regret bound of order  $\sqrt{KT}$  (and that the AdaHedge strategy does not need any mixing with the uniform distribution to achieve the bound of Theorem 2).

**Conclusion.** In this article, we considered adaptation to unknown ranges  $[m, M]$  in stochastic bandits. We did so for the sake of clarity but note that the real source of the difficulties stems from not knowing the upper bound  $M$  on the payoffs.

## 6. Numerical Experiments

We describe some numerical experiments on synthetic data to illustrate the performance of the algorithm(s) introduced compared to earlier approaches; we focus on how algorithms adapt to the scale of payoffs.

**Bandit problems considered and UCB strategies.** We consider stochastic bandit problems  $\underline{\nu}^{(\alpha)} = (\nu_a^{(\alpha)})_{a \in [K]}$  indexed by a scale parameter  $\alpha \in \{0.01, 0.1, 1, 10\}$ . We take  $K = 10$  arms, each arm  $a$  being associated with a truncated Gaussian distribution. Precisely, the distribution  $\nu_a^{(\alpha)}$  is the distribution of the variable

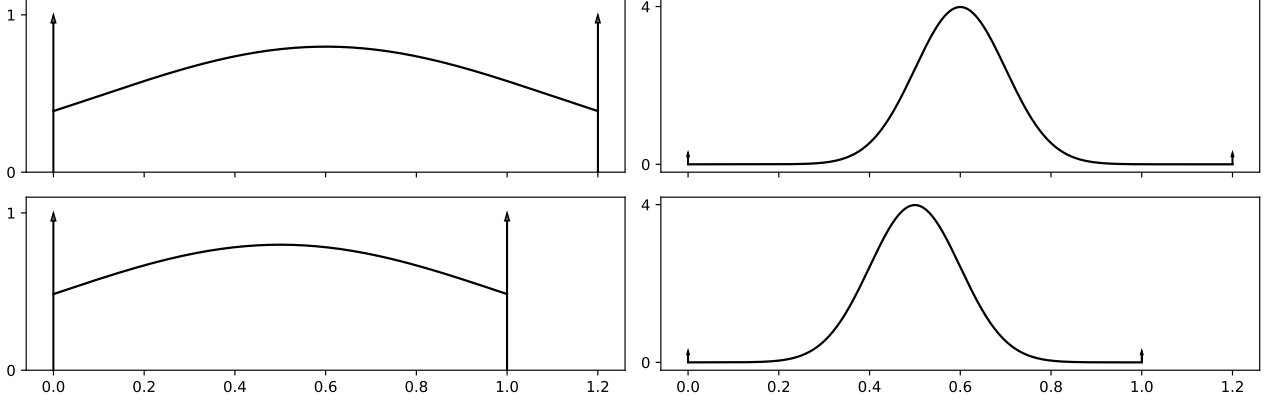
$$X_a^\alpha = \begin{cases} \alpha \max\{0, \min\{Y, 1.2\}\} & \text{with } Y \sim \mathcal{N}(0.6, V) \text{ if } a = 1, \\ \alpha \max\{0, \min\{Y, 1\}\} & \text{with } Y \sim \mathcal{N}(0.5, V) \text{ if } a \neq 1, \end{cases}$$

so that all distributions are commonly supported on  $[m, M] = [0, 1.2\alpha]$ , with arm 1 being the unique optimal arm. We will consider two values for  $V$ , namely  $V = 0.01$  (low-variance case) and  $V = 0.25$  (high-variance case). See Figure 1 for a plot of the corresponding probability density functions.

Given the lengths  $M - m = 1.2\alpha$  obtained for the ranges  $[m, M]$  as  $\alpha$  varies, we consider five instances of UCB (Auer et al., 2002a), with respective upper confidence bounds

$$\hat{\mu}_a(t) + 1.2\sigma\sqrt{\frac{2\ln T}{N_a(t)}}, \quad \text{for } \sigma \in \{0.001, 0.01, 0.1, 1, 10\},$$

where  $N_a(t)$  is the number of times arm  $a$  was pulled up to round  $t$  and  $\hat{\mu}_a(t)$  denotes the empirical average of payoffs obtained for arm  $a$  when it was played. Theory usually considers the tuning  $\sigma = \alpha$  (when  $\alpha$  is known).



**Figure 1:** Probability density functions of the reward distributions with respect to the sum of the Lebesgue measure and Dirac masses at 0, 1, and 1.2. Left pictures: high-variance case; right pictures: low-variance case. Top pictures: first arm (optimal arm); bottom pictures: other arms. Arrows represent atoms and their lengths are only illustrative.

**Range-estimating UCB** For the sake of completeness, we also plot the results of a version of UCB estimating the range, i.e., considering indices of the form

$$\hat{\mu}_a(t) + \hat{r}_t \sqrt{\frac{2 \ln T}{N_a(t)}}, \quad \text{where} \quad \hat{r}_t = \max_{s \leq t} Y_{A_s, s} - \min_{s \leq t} Y_{A_s, s}$$

estimates the range  $M - m$ . We were unable to provide theoretical guarantees that match our lower bounds, and this algorithm does not perform particularly well in practice as we will discuss below.

**Three other algorithms are considered.** For comparison, we also add a simple follow-the-leader strategy (referred to as FTL; i.e., a strategy picking at each round the arm with best payoff estimate so far) and the random strategy (i.e., picking at each round an arm uniformly at random). FTL and the random strategies will exhibit undesirable performance similar to the ones of incorrectly tuned instances of UCB (respectively, with too small and too large a parameter  $\sigma$ ).

The main algorithm of interest is, of course, the AHB strategy with extra-exploration (Algorithm 1), which we tune as indicated in Theorem 2 with parameter  $1/2$ .

**Experimental setting.** Each algorithm is run  $N = 300$  times, on a time horizon  $T = 100,000$ . We plot estimates of the rescaled regret  $R_T(\underline{\nu}^{(\alpha)})/\alpha$  to have a meaningful comparison between the bandit problems. These estimates are constructed as follows. We denote by  $\mu_1^{(\alpha)} = 0.6\alpha$  and  $\mu_a^{(\alpha)} = 0.5\alpha$  if  $a \neq 1$  the means associated with the distributions  $\nu_1^{(\alpha)}$  and  $\nu_a^{(\alpha)}$ , respectively. We index the arms picked in the  $n$ -th run by an additional subscript  $n$ , so that  $A_{T,n}$  refers to the arm picked by some strategy at time  $t$  in the  $n$ -th run. The expected regret of a given strategy can be rewritten as

$$R_T(\underline{\nu}^{(\alpha)}) = T \max_{a \in [K]} \mu_a^{(\alpha)} - \mathbb{E} \left[ \sum_{t=1}^T \mu_{A_t}^{(\alpha)} \right] = T \times (0.6\alpha) - \mathbb{E} \left[ \sum_{t=1}^T \mu_{A_t}^{(\alpha)} \right]$$

and is estimated by

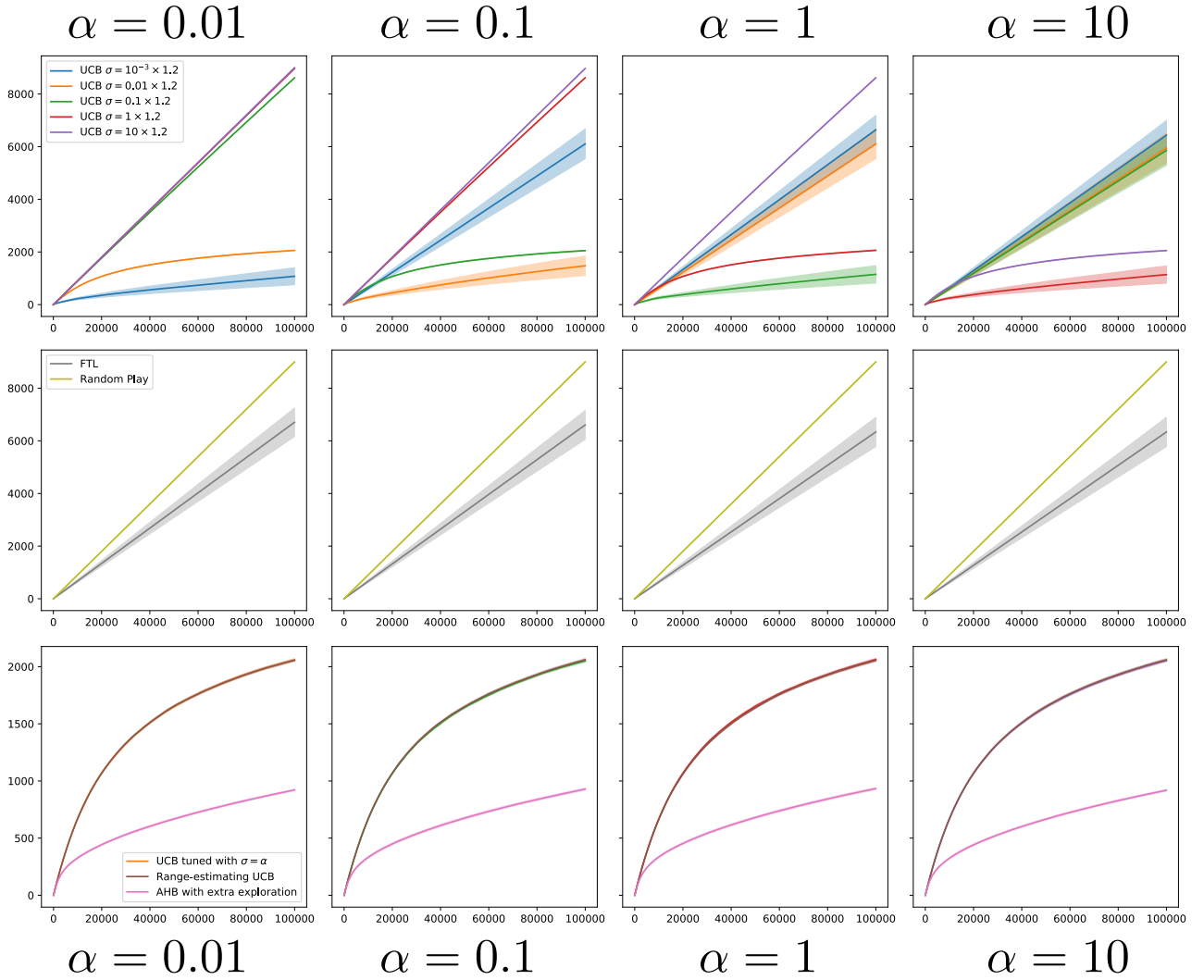
$$\hat{R}_T(\alpha) = \frac{1}{N} \sum_{n=1}^N \hat{R}_T(\alpha, n) \quad \text{where} \quad \hat{R}_T(\alpha, n) = T \times (0.6\alpha) - \sum_{t=1}^T \mu_{A_{t,n}}^{(\alpha)}.$$

On Figure 2 we plot the estimates  $\hat{R}_T(\alpha)/\alpha$  of the rescaled regret as solid lines. The shaded areas correspond to  $\pm 2$  standard errors of the sequences  $(\hat{R}_T(\alpha, n)/\alpha)_{n \in [N]}$ . All experiments were designed in Python, using the NumPy and joblib libraries, and were run on a standard laptop computer (with

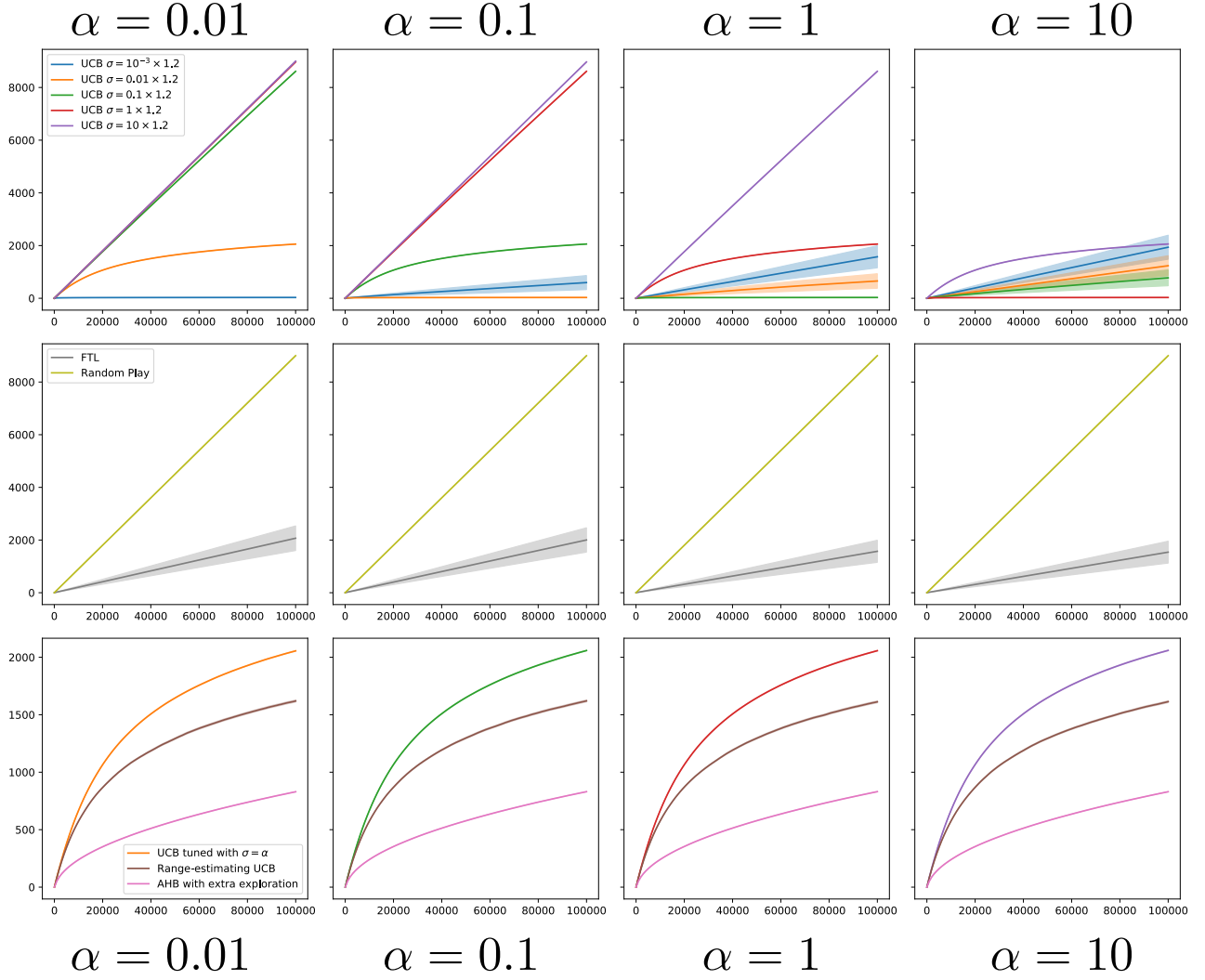
an Intel Core i5 processor). The code and setup for these experiments were only moderately optimized for computational efficiency. We display the average runtimes of all algorithms in Table 1; they are provided only for illustration and could certainly be significantly improved.

**Table 1:** Average runtimes of the (families of) algorithms considered, measured in seconds per run; as a reminder, we performed  $N = 300$  runs for each algorithm.

Random play	FTL	UCB	AHB
$X = 0.76 \text{ s/run}$	$1.8X$	$1.8X$	$6.5X$



**Figure 2:** Comparison of the (estimated) regrets of various strategies over bandit problems  $\mathcal{V}^{(\alpha)}$  in the high variance case, where  $\alpha$  ranges in  $\{0.01, 0.1, 1, 10\}$  and  $V = 0.25$ . Each algorithm was run  $N = 300$  times on every problem for  $T = 100,000$  time steps. Solid lines report the values of the estimated regrets, while shaded areas correspond to  $\pm 2$  standard errors of the estimates.



**Figure 3:** Same legend as for Figure 2, but in the low variance case.

**Discussion of the results.** A first observation is that, as expected, AHB (see the third lines of Figures 2 and 3) is unaffected by the scale of the problems. The same can be said for FTL, the random strategy, and the range-estimating UCB.

AHB yields favorable results (note that the range of the  $y$ -axis for the third line is smaller than the ranges in the first two lines), exhibiting better results in all situations than UCB tuned with the correct scale  $\sigma = \alpha$  and the range-estimating UCB. Note that the latter two strategies behave in a similar manner in the low-variance case and have virtually indistinguishable performance in the high-variance case.

Our major second observation is that the performance of UCB depends dramatically on the value of the parameter  $\sigma$ . When the scale parameter  $\sigma$  is too large, UCB is essentially playing at random, i.e., overexploring, and incurs the same linear regret as the random strategy. On the other hand, when UCB is run with too small a scale parameter  $\sigma$ , it underexplores and incurs a large regret, which we discuss in greater details. We measure the cost of underexploration through the performance of the instances of UCB tuned with  $\sigma = \alpha/10$  (moderate underexploration) and  $\sigma = \alpha/100$  (severe underexploration). In the the high-variance case, (moderate or severe) underexploration leads to linear regret; see the first line of Figure 2. In the low variance case, moderate underexploration is beneficial while severe underexploration leads again to linear regret.

## A. Proofs of Two Claims on Distribution-Dependent Regret Rates

In this section, we expand on the two claims stated after Definition 2, that are relative to distribution-dependent lower bounds for adaptation to the range: first, that all reasonable strategies (in the sense of Definition 4 below with  $\mathcal{D} = \mathcal{D}_{-,+}$ ) ensure that for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}_{-,+}$  with at least one suboptimal arm,

$$\liminf_{T \rightarrow +\infty} \frac{R_T(\underline{\nu})}{\ln T} = +\infty,$$

while, second, any rate  $\Phi_{\text{dep}} \gg \ln T$  may be achieved thanks to a simple upper-confidence bound [UCB] strategy. Before we do so, we remind the reader of the “classical” results, for an abstract model  $\mathcal{D}$  and then, for the model  $\mathcal{D}_{m,M}$  corresponding to payoff distributions with a known range  $[m, M]$ .

**Definition 4.** *A strategy is uniformly fast convergent on a model  $\mathcal{D}$  if for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}$ , it achieves a subpolynomial regret bound, that is,  $R_T(\underline{\nu})/T^\alpha \rightarrow 0$  for all  $(\alpha, 1]$ .*

A lower bound on the distribution-dependent rates that such a strategy may achieve is provided by a general result of Lai and Robbins [1985] and Burnetas and Katehakis [1996]; see also its rederivation by Garivier et al. [2019b]. It involves a quantity defined as an infimum of Kullback-Leibler divergences: we recall that for two probability distributions  $\nu, \nu'$  defined on the same probability space  $(\Omega, \mathcal{F})$ ,

$$\text{KL}(\nu, \nu') = \begin{cases} \int_{\Omega} \ln\left(\frac{d\nu}{d\nu'}\right) d\nu & \text{if } \nu \ll \nu', \\ +\infty & \text{otherwise,} \end{cases}$$

where  $\nu \ll \nu'$  means that  $\nu$  is absolutely continuous with respect to  $\nu'$  and  $d\nu/d\nu'$  then denotes the Radon-Nikodym derivative. Now, for any probability distribution  $\nu$ , any real number  $x$ , and any model  $\mathcal{D}$ , we define

$$\mathcal{K}_{\text{inf}}(\nu, x, \mathcal{D}) = \inf\{\text{KL}(\nu, \nu') : \nu' \in \mathcal{D} \text{ and } \mathbb{E}(\nu') > x\},$$

where by convention, the infimum of an empty set equals  $+\infty$  and where we denoted by  $\mathbb{E}(\nu')$  the expectation of  $\nu'$ . The quantity  $\mathcal{K}_{\text{inf}}(\nu, x, \mathcal{D})$  can be null. With the usual measure-theoretic conventions, in particular,  $0/0 = 0$ , we then have the following lower bound.

**Reminder 1.** *For all models  $\mathcal{D}$ , for all uniformly fast convergent strategies on  $\mathcal{D}$ , for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}$ ,*

$$\liminf_{T \rightarrow +\infty} \frac{R_T(\underline{\nu})}{\ln T} \geq \sum_{a \in [K]} \frac{\Delta_a}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*, \mathcal{D})}.$$

When the range  $[m, M]$  of payoffs is known, i.e., when the model is  $\mathcal{D}_{m,M}$ , there exist strategies achieving the lower bound of Reminder 1, like the DMED strategy of Honda and Takemura [2011, 2015] or the KL-UCB strategy of Cappé et al. [2013] and Garivier et al. [2019a]. (This actually even holds for semi-bounded payoffs with a known upper bound on the payoffs, as is discussed in details in an online appendix [arXiv:2006.03378], Section E.1.)

### A.1. No Logarithmic Distribution-Dependent Regret Bound under Adaptation to the Range

Now, the lower bound in Reminder 1 equals  $+\infty$  when the range is not known, that is, when we consider the model  $\mathcal{D}_{-,+}$  of bounded distributions with unknown range. Actually, the proof reveals (just like the proof of Theorem 1 in Section 3) that the important fact is that the upper end of the payoff range is unknown. The impossibility result also holds for models  $\mathcal{D}_{m,+}$  of bounded distributions with unknown upper end on the range but known lower end  $m$  on the range, which were defined in (7).



**Theorem 4.** *All uniformly fast convergent strategies on  $\mathcal{D}_{-,+}$  are such that, for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}_{-,+}$  with at least one suboptimal arm  $a$ ,*

$$\mathcal{K}_{\inf}(\nu_a, \mu^*, \mathcal{D}_{m,+}) = 0, \quad \text{thus} \quad \liminf_{T \rightarrow +\infty} \frac{R_T(\underline{\nu})}{\ln T} = +\infty.$$

*The same result holds for all models  $\mathcal{D}_{m,+}$ , where  $m \in \mathbb{R}$ .*

Strategies that are adaptive to the range thus cannot get rates  $\Phi_{\text{dep}}$  for distribution-dependent regret bounds of the order of  $\ln T$ . A similar phenomenon was discussed by Lattimore [2017] in the case of stochastic bandits with sub-Gaussian distributions. It however turns out that any rate  $\Phi_{\text{dep}}$  such that  $\Phi_{\text{dep}}(T) \gg \ln T$  may be achieved, through a simple upper-confidence bound [UCB] strategy, as also discussed by Lattimore [2017]; see further details after the proof.

Interestingly, Cowan and Katehakis [2015] observe that for the model of uniform distributions over intervals, the  $\mathcal{K}_{\inf}$  is positive, and thus the lower bound of Reminder 1 does not prevent logarithmic regret bounds. In fact, they also provide an algorithm enjoying optimal distribution-dependent bounds—thus being, in a sense, adaptive to the range in that very restricted model.

*Proof.* We fix  $m \in \mathbb{R}$  and provide the proof for  $\mathcal{D}_{m,+}$ . We show below that  $\mathcal{K}_{\inf}(\nu_a, \mu^*, \mathcal{D}_{m,+}) = 0$  for any suboptimal  $a$  with  $\nu_a \in \mathcal{D}_{m,+}$ . The  $\liminf$  being equal to  $+\infty$  follows from Reminder 1 and the fact  $\Delta_a = \mu^* - \mu_a > 0$  as  $a$  is suboptimal.

We have in particular  $\mu_a \geq m$ . We use the same construction as in the proof of Theorem 1. Let  $\nu'_\varepsilon = (1 - \varepsilon)\nu_a + \varepsilon\delta_{\mu_a + 2\Delta_a/\varepsilon}$  for  $\varepsilon \in (0, 1)$ : it is a bounded probability distribution, with lower end of support larger than  $m$ , that is,  $\nu'_\varepsilon \in \mathcal{D}_{m,+}$ . For  $\varepsilon$  small enough,  $\mu_a + 2\Delta_a/\varepsilon$  lies outside of the bounded support of  $\nu_a$ . In that case, the density of  $\nu_a$  with respect to  $\nu'_\varepsilon$  is given by  $1/(1 - \varepsilon)$  on the support of  $\nu_a$  (and 0 elsewhere), so that

$$\text{KL}(\nu_a, \nu'_\varepsilon) = \ln\left(\frac{1}{1 - \varepsilon}\right).$$

Moreover,  $\mathbb{E}(\nu'_\varepsilon) = (1 - \varepsilon)\mu_a + \varepsilon(\mu_a + 2\Delta_a/\varepsilon) = \mu_a + 2\Delta_a = \mu^* + \Delta_a > \mu^*$ . Therefore, by definition of  $\mathcal{K}_{\inf}$  as an infimum,

$$\mathcal{K}_{\inf}(\nu_a, \mu^*, \mathcal{D}_{m,+}) \leq \text{KL}(\nu_a, \nu'_\varepsilon) = \ln\left(\frac{1}{1 - \varepsilon}\right).$$

This upper bound holds for all  $\varepsilon > 0$  small enough and thus shows that  $\mathcal{K}_{\inf}(\nu_a, \mu^*, \mathcal{D}_{m,+}) = 0$ .

The exact same construction and proof can be performed in the case of  $\mathcal{D}_{-,+}$ , without the need of indicating that the lower end of the support of  $\nu'_\varepsilon$  is larger than  $m$ .  $\square$

## A.2. UCB with an Increased Exploration Rate Adapts to the Range

The lower bound of Theorem 4 does not prevent distribution-dependent rates for adaptation that are arbitrarily larger than a logarithm. Consider UCB with indexes of the form

$$\hat{\mu}_a(t) + \sqrt{\frac{\varphi(t)}{N_a(t)}} \quad \text{where} \quad \frac{\varphi(t)}{\ln t} \rightarrow +\infty \quad \text{and} \quad \frac{\varphi(t)}{t} \rightarrow 0,$$

and where  $\hat{\mu}_a(t)$  denotes the empirical average of payoffs obtained till round  $t$  when playing arm  $a$ . Following the analysis of Lattimore [2017] in the case of Gaussian bandits with unknown variances, it can be shown that such a UCB is adaptive to the unknown range of payoffs with a distribution-dependent rate  $\Phi_{\text{dep}} = \varphi$ . However the trick used here is purely asymptotic and gives up on finite-time guarantees.

## B. Proof of Theorem 2

**How the second regret bound follows from the first one.** We substitute the indicated values of the  $\gamma_t$ . We have, first,

$$\sum_{t=K+1}^T \gamma_t \leq \sqrt{5(1-\alpha)K \ln K} \sum_{t=K+1}^T t^{-\alpha} \leq \sqrt{5(1-\alpha)K \ln K} \int_0^T \frac{1}{t^\alpha} dt = \sqrt{\frac{5K \ln K}{1-\alpha}} T^{1-\alpha}, \quad (14)$$

second, using the definition of  $\gamma_T$  as a minimum,

$$\frac{K \ln K}{\gamma_T} \leq \frac{K \ln K}{1/2} + \frac{T^\alpha K \ln K}{\sqrt{5(1-\alpha)K \ln K}} = 2K \ln K + \sqrt{\frac{K \ln K}{5(1-\alpha)}} T^\alpha,$$

and third,  $\sqrt{T} \leq T^{\max\{\alpha, 1-\alpha\}}$ , so that the first regret bound of Theorem 2 is further bounded by

$$(M-m)\sqrt{K \ln K} \left( 3 + 2\sqrt{\frac{5}{1-\alpha}} \right) T^{\max\{\alpha, 1-\alpha\}} + 10(M-m)K \ln K.$$

The claimed expression for  $\Phi_{\text{adv}}(T)$  is obtained by bounding  $2\sqrt{5}$  by 5.

**First regret bound.** In Algorithm 1, for time steps  $t \geq K+1$ , the weights  $q_t$  are obtained by using the AdaHedge algorithm of De Rooij et al. [2014] on the payoff estimates  $\hat{y}_{t,a}$ . AdaHedge is designed for the case of a full monitoring (not a bandit monitoring), but the use of these estimates emulates a full monitoring. Section 2.2 of De Rooij et al. [2014]—see also an earlier analysis by Cesa-Bianchi et al. [2007]—ensures the bound stated next in Reminder 2. (For the sake of completeness, we rederive this bound in an online appendix [arXiv:2006.03378], Section E.2.2.) We call pre-regret the quantity at hand in Reminder 2: it corresponds to some regret defined in terms of the payoff estimates.

**Reminder 2** (Application of Lemma 3 and Theorem 6 of De Rooij et al., 2014). *For all sequences of payoff estimates  $\hat{y}_{t,a}$  lying in some bounded real-valued interval, denoted by  $[b, B]$ , for all  $T \geq K+1$ , the pre-regret of AdaHedge satisfies*

$$\begin{aligned} \max_{k \in [K]} \sum_{t=K+1}^T \hat{y}_{t,k} - \sum_{t=K+1}^T \sum_{a=1}^K q_{t,a} \hat{y}_{t,a} &\leq 2 \sum_{t=K+1}^T \delta_t \\ \text{where} \quad \sum_{t=K+1}^T \delta_t &\leq \underbrace{\sqrt{\sum_{t=K+1}^T \sum_{a=1}^K q_{t,a} \left( \hat{y}_{t,a} - \sum_{k \in [K]} q_{t,k} \hat{y}_{t,k} \right)^2 \ln K}}_{\leq \sqrt{\sum_{t=K+1}^T \sum_{a=1}^K q_{t,a} (\hat{y}_{t,a} - c)^2 \ln K} \quad \text{for any } c \in \mathbb{R}} + (B-b) \left( 1 + \frac{2}{3} \ln K \right) \end{aligned}$$

and AdaHedge does not require the knowledge of  $[b, B]$  to achieve this bound.

The bound of Reminder 2 will prove itself particularly handy for three reasons: first, it is valid for signed payoffs (payoffs in  $\mathbb{R}$ ); second, it is adaptive to the range of payoffs; third, the right-hand side looks at first sight not intrinsic enough a bound (as it also depends on the weights  $q_t$ ) but we will see later that this dependency is particularly useful.

We recall that we start the summation in Reminder 2 at  $t = K+1$  because the AdaHedge algorithm is only started at this time, after the initial exploration. The bound holding “for any  $c \in \mathbb{R}$ ” is obtained by a classical bound on the variance.

*Proof of the first bound of Theorem 2.* We deal with the contribution of the initial exploration by using the inequality  $\max(u + v) \leq \max u + \max v$ , together with the fact that  $y_{t,a} - y_{t,A_t} \leq M - m$  for any  $a \in [K]$ :

$$R_T(y_{1:T}) \leq \underbrace{\max_{a \in [K]} \sum_{t=1}^K y_{t,a} - \mathbb{E} \left[ \sum_{t=1}^K y_{t,A_t} \right]}_{\leq K(M-m)} + \max_{a \in [K]} \sum_{t=K+1}^T y_{t,a} - \mathbb{E} \left[ \sum_{t=K+1}^T y_{t,A_t} \right]. \quad (15)$$

We now transform the pre-regret bound of Reminder 2, which is stated with the distributions  $q_t$ , into a pre-regret bound with the distributions  $p_t$ ; we do so while substituting the bounds  $B = C + KM/\gamma_T$  and  $b = C + Km/\gamma_T$  implied by (11) and the fact that  $(\gamma_t)$  is non-increasing, and by using the definition  $q_{t,a} = p_{t,a} - \gamma_t(1/K - q_{t,a})$  for all  $a \in [K]$ :

$$\begin{aligned} \max_{k \in [K]} \sum_{t=K+1}^T \hat{y}_{t,k} - \sum_{t=K+1}^T \sum_{a=1}^K p_{t,a} \hat{y}_{t,a} + \sum_{t=K+1}^T \gamma_t \sum_{a=1}^K (1/K - q_{t,a}) \hat{y}_{t,a} &\leq 2 \sum_{t=K+1}^T \delta_t \\ \text{where } \sum_{t=K+1}^T \delta_t &\leq \sqrt{\sum_{t=K+1}^T \sum_{a=1}^K q_{t,a} (\hat{y}_{t,a} - C)^2 \ln K} + \frac{(M-m)K}{\gamma_T} \left( 1 + \frac{2}{3} \ln K \right). \end{aligned} \quad (16)$$

As noted by Auer et al. [2002b], by the very definition (9) of the estimates,

$$\sum_{a=1}^K p_{t,a} \hat{y}_{t,a} = y_{t,A_t}.$$

By (10), the tower rule and the fact that  $q_t$  is  $H_{t-1}$ -measurable, on the one hand, and the fact that the expectation of a maximum is larger than the maximum of expectations, on the other hand, the left-hand side of the first inequality in (16) thus satisfies

$$\begin{aligned} &\mathbb{E} \left[ \max_{k \in [K]} \sum_{t=K+1}^T \hat{y}_{t,k} - \sum_{t=K+1}^T \sum_{a=1}^K p_{t,a} \hat{y}_{t,a} + \sum_{t=K+1}^T \gamma_t \sum_{a=1}^K (1/K - q_{t,a}) \hat{y}_{t,a} \right] \\ &\geq \max_{k \in [K]} \sum_{t=K+1}^T y_{t,k} - \mathbb{E} \left[ \sum_{t=K+1}^T y_{t,A_t} \right] + \sum_{t=K+1}^T \gamma_t \left( \underbrace{\sum_{a=1}^K y_{t,a}/K}_{\in [m,M]} - \underbrace{\sum_{a=1}^K \mathbb{E}[q_{t,a}] y_{t,a}}_{\in [m,M]} \right) \\ &\geq \max_{k \in [K]} \sum_{t=K+1}^T y_{t,k} - \mathbb{E} \left[ \sum_{t=K+1}^T y_{t,A_t} \right] - (M-m) \sum_{t=1}^T \gamma_t. \end{aligned}$$

As for the right-hand side of the second inequality in (16), we first note that by definition (see line 4 in Algorithm 1),  $p_{t,a} \geq (1 - \gamma_t)q_{t,a}$  with  $\gamma_t \leq 1/2$  by assumption on the extra-exploration rate, so that  $q_{t,a} \leq 2p_{t,a}$ ; therefore, by substituting first this inequality and then by using Jensen's inequality,

$$\begin{aligned} \mathbb{E} \left[ \sqrt{\sum_{t=K+1}^T \sum_{a=1}^K q_{t,a} (\hat{y}_{t,a} - C)^2 \ln K} \right] &\leq \sqrt{2} \mathbb{E} \left[ \sqrt{\sum_{t=K+1}^T \sum_{a=1}^K p_{t,a} (\hat{y}_{t,a} - C)^2 \ln K} \right] \\ &\leq \sqrt{2} \sqrt{\sum_{t=K+1}^T \sum_{a=1}^K \mathbb{E} [p_{t,a} (\hat{y}_{t,a} - C)^2] \ln K}. \end{aligned} \quad (17)$$

Standard calculations (see Auer et al., 2002b again) show, similarly to (10), that for all  $a \in [K]$ ,

$$\mathbb{E} \left[ p_{t,a} (\hat{y}_{t,a} - C)^2 \mid H_{t-1} \right] = \mathbb{E} \left[ \frac{(y_{t,A_t} - C)^2}{p_{t,a}} \mathbf{1}_{\{A_t=a\}} \right] = (y_{t,a} - C)^2 \leq (M-m)^2,$$

where the last inequality comes from (11). By the tower rule, the same upper bound holds for the (unconditional) expectation. Therefore, taking the expectation of both sides of (16) and collecting all bounds together, we proved so far

$$R_T(y_{1:T}) \leq \underbrace{2\sqrt{2}}_{\leq 3} (M - m) \sqrt{KT \ln K} + (M - m) \frac{K \ln K}{\gamma_T} \underbrace{\left( \frac{2 + \gamma_T}{\ln K} + \frac{4}{3} \right)}_{\leq 5} + (M - m) \sum_{t=K+1}^T \gamma_t,$$

where we used  $\gamma_T \leq 1/2$  and  $\ln K \geq \ln 2$  as  $K \geq 2$ .  $\square$

### C. Proof of Theorem 3

Given the decomposition (1) of the regret, it is necessary and sufficient to upper bound the expected number of times  $\mathbb{E}[N_a(t)]$  any suboptimal arm  $a$  is drawn, where by definition of Algorithm 1,

$$\mathbb{E}[N_a(t)] = 1 + \mathbb{E} \left[ \sum_{t=K+1}^T \left( (1 - \gamma_t) q_{t,a} + \frac{\gamma_t}{K} \right) \right] \leq 1 + \sum_{t=K+1}^T \mathbb{E}[q_{t,a}] + \frac{1}{K} \sum_{t=K+1}^T \gamma_t.$$

We show below (and this is the main part of the proof) that

$$\sum_{t=K+1}^T \mathbb{E}[q_{t,a}] = \mathcal{O}(\ln T). \quad (18)$$

The straightforward calculations (14) already showed that

$$\frac{1}{K} \sum_{t=K+1}^T \gamma_t \leq \sqrt{\frac{5 \ln K}{(1 - \alpha)K}} T^{1-\alpha}.$$

Substituting the value (12) of  $\Phi_{\text{free}}^{\text{AHB}}(T) = \Phi_{\text{adv}}(T)$  and using the decomposition (1) of  $R_T(\underline{\nu})$  into  $\sum \Delta_a \mathbb{E}[N_a(t)]$  then yield

$$\frac{R_T(\underline{\nu})}{T/\Phi_{\text{free}}^{\text{AHB}}(T)} \leq \sum_{a \in [K]} \Delta_a \sqrt{\frac{5 \ln K}{(1 - \alpha)K}} \left( 3 + \frac{5}{\sqrt{1 - \alpha}} \right) \sqrt{K \ln K} (1 + o(1)) + \mathcal{O} \left( \frac{\ln T}{T^{1-\alpha}} \right),$$

from which the stated bound follows, via the crude inequality  $3\sqrt{5}\sqrt{1 - \alpha} + 5 \leq 12$ .

**Structure of the proof of (18).** Let  $a^*$  denote an optimal arm. By definition of  $q_{t,a}$  and by lower bounding a sum of exponential terms by any of the summands, we get

$$q_{t,a} = \frac{\exp \left( \eta_t \sum_{s=K+1}^{t-1} \hat{y}_{t,a} \right)}{\sum_{k=1}^K \exp \left( \eta_t \sum_{s=K+1}^{t-1} \hat{y}_{t,k} \right)} \leq \exp \left( \eta_t \sum_{t=K+1}^{t-1} (\hat{y}_{t,a} - \hat{y}_{t,a^*}) \right).$$

Then, by separating cases, depending on whether  $\sum_{t=K+1}^{t-1} (\hat{y}_{t,a} - \hat{y}_{t,a^*})$  is smaller or larger than the threshold  $-(t - 1 - K)\Delta_a/2$ , and by remembering that the probability  $q_{t,a}$  is always smaller than 1, we get

$$\begin{aligned} \sum_{t=K+1}^T \mathbb{E}[q_{t,a}] &\leq \sum_{t=K+1}^T \mathbb{E} \left[ \exp \left( -\eta_t \frac{(t - 1 - K)\Delta_a}{2} \right) \right] \\ &\quad + \sum_{t=K+1}^T \mathbb{P} \left[ \sum_{s=K+1}^{t-1} (\hat{y}_{s,a} - \hat{y}_{s,a^*}) \geq -\frac{(t - 1 - K)\Delta_a}{2} \right]. \end{aligned} \quad (19)$$

We show that the sums in the right-hand side of (19) are respectively  $\mathcal{O}(1)$  and  $\mathcal{O}(\ln T)$ .

**First sum in the right-hand side of (19).** Given the definition of the learning rates (see the statement of Algorithm 1), namely,

$$\eta_t = \ln K \left/ \sum_{s=K+1}^{t-1} \delta_s \right., \quad (20)$$

we are interested in upper bounds on the sum of the  $\delta_s$ . Such upper bounds were already derived in the proof of Theorem 2; the second inequality in (16) together with the bound  $q_{t,a} \leq 2p_{t,a}$  stated in the middle of the proof immediately yield

$$\begin{aligned} \sum_{s=K+1}^{t-1} \delta_s &\leq \sqrt{\sum_{s=K+1}^t \sum_{a=1}^K q_{s,a} (\hat{y}_{s,a} - C)^2 \ln K} + \frac{(M-m)K}{\gamma_t} \left(1 + \frac{2}{3} \ln K\right) \\ &\leq \sqrt{2} \sqrt{\sum_{s=K+1}^t \sum_{a=1}^K p_{s,a} (\hat{y}_{s,a} - C)^2 \ln K} + \frac{(M-m)K}{\gamma_t} \left(1 + \frac{2}{3} \ln K\right). \end{aligned}$$

Unlike what we did to complete the proof of Theorem 2, we do not take expectations and rather proceed with deterministic bounds. By the definition (9) of the estimated payoffs for the equality below, by (11) for the first inequality below, and by the fact that the exploration rates are non-increasing for the second inequality below, we have, for all  $s \geq K+1$ ,

$$\sum_{a=1}^K p_{s,a} (\hat{y}_{s,a} - C)^2 = \frac{(y_{s,A_s} - C)^2}{p_{s,A_s}} \leq \frac{(M-m)^2}{\gamma_s/K} \leq \frac{(M-m)^2}{\gamma_t/K}. \quad (21)$$

Therefore,

$$\sum_{s=K+1}^{t-1} \delta_s \leq \sqrt{2}(M-m) \sqrt{\frac{tK \ln K}{\gamma_t}} + \frac{(M-m)K}{\gamma_t} \left(1 + \frac{2}{3} \ln K\right) \stackrel{\text{def}}{=} D_t = \Theta\left(\sqrt{t/\gamma_t} + 1/\gamma_t\right).$$

For the sake of concision, we denoted by  $D_t$  the obtained bound. Via the definition (20) of  $\eta_t$ , the sum of interest is in turn bounded by

$$\sum_{t=K+1}^T \exp\left(-\eta_t(t-1-K) \frac{\Delta_a}{2}\right) \leq \sum_{t=K+1}^T \exp\left(-\frac{\Delta_a \ln K}{2} \frac{t-1-K}{D_t}\right) = \mathcal{O}(1),$$

where the equality to  $\mathcal{O}(1)$ , i.e., the fact that the considered series is bounded, follows from the fact that

$$-(t-1-K)/D_t = \Theta\left(\sqrt{t\gamma_t} + t\gamma_t\right) = \Theta(t^{(1-\alpha)/2} + t^{1-\alpha}).$$

**Second sum in the right-hand side of (19).** We will use Bernstein's inequality for martingales, and more specifically, the formulation of the inequality by Freedman [1975, Thm. 1.6]—see also Massart [2007, Section 2.2]—, as stated next.

**Reminder 3.** Let  $(X_n)_{n \geq 1}$  be a martingale difference sequence with respect to a filtration  $(\mathcal{F}_n)_{n \geq 0}$ , and let  $N \geq 1$  be a summation horizon. Assume that there exist real numbers  $b$  and  $v_N$  such that, almost surely,

$$\forall n \leq N, \quad X_n \leq b \quad \text{and} \quad \sum_{n=1}^N \mathbb{E}[X_n^2 | \mathcal{F}_{n-1}] \leq v_N.$$

Then for all  $\delta \in (0, 1)$ ,

$$\mathbb{P}\left[\sum_{n=1}^N X_n \geq \sqrt{2v_N \ln \frac{1}{\delta}} + \frac{b}{3} \ln \frac{1}{\delta}\right] \leq \delta.$$

For  $s \geq K + 1$ , we consider the increments  $X_s = \Delta_a - \hat{y}_{s,a^*} + \hat{y}_{s,a}$ , which are adapted to the filtration  $\mathcal{F}_s = \sigma(A_1, Z_1, \dots, A_s, Z_s)$ , where we recall that  $Z_1, \dots, Z_s$  denote the payoffs obtained in rounds  $1, \dots, s$ . Also, as  $p_s$  is measurable with respect to past information  $\mathcal{F}_{s-1}$  and since payoffs are drawn independently from everything else (see Section 2), we have, by the definition (9) of the estimated payoffs (where we rather denote by  $Y_{s,a}$  the payoffs drawn at random according to  $\nu_a$ , to be in line with the notation of Section 2 for stochastic bandits): for all  $a \in [K]$ ,

$$\mathbb{E}[\hat{y}_{s,a} \mid \mathcal{F}_{s-1}] = \frac{\mathbb{E}[Y_{s,a} \mid \mathcal{F}_{s-1}] - C}{p_{s,a}} \mathbb{1}_{\{A_s=a\}} + C = \frac{\mu_a - C}{p_{s,a}} \mathbb{1}_{\{A_s=a\}} + C = \mu_a.$$

As a consequence,  $\mathbb{E}[X_s \mid \mathcal{F}_{s-1}] = \mathbb{E}[\Delta_a - \hat{y}_{s,a^*} + \hat{y}_{s,a} \mid \mathcal{F}_{s-1}] = 0$ . Put differently,  $(X_s)_{s \geq K+1}$  is indeed a martingale difference sequence with respect to the filtration  $(\mathcal{F}_s)_{s \geq K}$ .

We now check that the additional assumptions of Reminder 3 are satisfied. Manipulations and arguments similar to the ones used in (11) and (21) show that for all  $s \geq K + 1$ ,

$$\begin{aligned} \Delta_a - \hat{y}_{s,a^*} + \hat{y}_{s,a} &\leq \Delta_a - \frac{Y_{s,a^*} - C}{p_{s,a}} \mathbb{1}_{\{A_s=a^*\}} + \frac{Y_{s,a} - C}{p_{s,a}} \mathbb{1}_{\{A_s=a\}} \\ &\leq (M - m)(1 + K/\gamma_s) \leq b \stackrel{\text{def}}{=} (M - m)(1 + K/\gamma_t). \end{aligned}$$

For the variance bound, we first note that for all  $s \leq t - 1$ , we have  $(\hat{y}_{s,a} - C)(\hat{y}_{s,a^*} - C) = 0$  because of the indicator functions, and therefore,

$$\begin{aligned} \mathbb{E}[(\Delta_a - \hat{y}_{s,a^*} + \hat{y}_{s,a})^2 \mid \mathcal{F}_{s-1}] &\leq \mathbb{E}[(\hat{y}_{s,a^*} + \hat{y}_{s,a})^2 \mid \mathcal{F}_{s-1}] \\ &\leq \mathbb{E}[(\hat{y}_{s,a^*} - C)^2 \mid \mathcal{F}_{s-1}] + \mathbb{E}[(\hat{y}_{s,a} - C)^2 \mid \mathcal{F}_{s-1}]; \end{aligned}$$

in addition, for all  $a \in [K]$  (including  $a^*$ ),

$$\mathbb{E}[(\hat{y}_{s,a} - C)^2 \mid \mathcal{F}_{s-1}] = \mathbb{E}\left[\frac{(Y_{s,A_s} - C)^2}{p_{s,a}^2} \mathbb{1}_{\{A_s=a\}} \mid \mathcal{F}_{s-1}\right] \leq \frac{(M - m)^2}{p_{s,a}} \leq \frac{(M - m)^2 K}{\gamma_t}.$$

Therefore

$$\sum_{s=K+1}^{t-1} \mathbb{E}[(\Delta_a - \hat{y}_{s,a^*} + \hat{y}_{s,a})^2 \mid \mathcal{F}_{s-1}] \leq \frac{2K(M - m)^2(t - 1 - K)}{\gamma_t} \leq v_t \stackrel{\text{def}}{=} \frac{2(M - m)^2 t K}{\gamma_t}.$$

Bernstein's inequality (Reminder 3) may thus be applied; the choice  $\delta = 1/t$  therein leads to

$$\mathbb{P}\left[\sum_{s=K+1}^{t-1} (\Delta_a - (\hat{y}_{s,a^*} - \hat{y}_{s,a})) \geq \underbrace{2(M - m) \sqrt{\frac{tK}{\gamma_t} \ln t} + \frac{M - m}{3} \left(1 + \frac{K}{\gamma_t}\right) \ln t}_{\stackrel{\text{def}}{=} D'_t}\right] \leq \frac{1}{t}.$$

As  $\sqrt{t/\gamma_t} = \mathcal{O}(t^{(1+\alpha)/2})$  and  $1/\gamma_t = \mathcal{O}(t^\alpha)$  as  $t \rightarrow \infty$ , where  $\alpha < 1$ , and as  $\Delta_a > 0$  (given that we are considering a suboptimal arm  $a$ ), there exists  $t_0 \in \mathbb{N}$  such that for all  $t \geq t_0$ ,

$$D'_t \leq \frac{(t - 1 - K)\Delta_a}{2}$$

thus

$$\begin{aligned} \mathbb{P}\left[\sum_{s=K+1}^{t-1} (\hat{y}_{s,a} - \hat{y}_{s,a^*}) \geq -\frac{(t - 1 - K)\Delta_a}{2}\right] &= \mathbb{P}\left[\sum_{s=K+1}^{t-1} (\Delta_a - (\hat{y}_{s,a^*} - \hat{y}_{s,a})) \geq \frac{(t - 1 - K)\Delta_a}{2}\right] \\ &\leq \mathbb{P}\left[\sum_{s=K+1}^{t-1} (\Delta_a - (\hat{y}_{s,a^*} - \hat{y}_{s,a})) \geq D'_t\right] \leq \frac{1}{t}. \end{aligned}$$

Therefore, as  $T \rightarrow \infty$

$$\sum_{t=1}^T \mathbb{P} \left[ \sum_{t=K+1}^{t-1} (\hat{y}_{t,a} - \hat{y}_{t,a^*}) \geq -\frac{(t-1-K)\Delta_a}{2} \right] = \mathcal{O}(\ln T),$$

as claimed. This concludes the proof.

## References

- J.-Y. Audibert. Fast learning rates in statistical inference through aggregation. *Annals of Statistics*, 37(4):1591–1646, 2009.
- J.-Y. Audibert and S. Bubeck. Minimax policies for adversarial and stochastic bandits. In *Proceedings of the 22nd Annual Conference on Learning Theory (COLT’09)*, pages 217–226. Omnipress, 2009.
- J.-Y. Audibert, R. Munos, and C. Szepesvári. Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):1876–1902, 2009.
- J.-Y. Audibert, S. Bubeck, and G. Lugosi. Regret in online combinatorial optimization. *Mathematics of Operations Research*, 39(1):31–45, 2014.
- P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002a.
- P. Auer, N. Cesa-Bianchi, Y. Freund, and R.E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002b.
- S. Bubeck, M.B. Cohen, and Y. Li. Sparsity, variance and curvature in multi-armed bandits. In *Proceedings of the 29th International Conference on Algorithmic Learning Theory (ALT’18)*, volume 83 of PMLR, pages 111–127, 2018.
- A.N. Burnetas and M.N. Katehakis. Optimal adaptive policies for sequential allocation problems. *Advances in Applied Mathematics*, 17(2):122–142, 1996.
- O. Cappé, A. Garivier, O.-A. Maillard, R. Munos, and G. Stoltz. Kullback–Leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, 41(3):1516–1541, 2013.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- N. Cesa-Bianchi and O. Shamir. Bandit regret scaling with the effective loss range. In *Proceedings of the 29th International Conference on Algorithmic Learning Theory (ALT’18)*, volume 83 of PMLR, pages 128–151, 2018.
- N. Cesa-Bianchi, Y. Mansour, and G. Stoltz. Improved second-order bounds for prediction with expert advice. *Machine Learning*, 66(2-3):321–352, 2007.
- Y. Chow and H. Teicher. *Probability Theory*. Springer, 1988.
- W. Cowan and M.N. Katehakis. An asymptotically optimal policy for uniform bandits of unknown support, 2015. Preprint, arXiv:1505.01918.
- W. Cowan, J. Honda, and M.N. Katehakis. Normal bandits of unknown means and variances. *Journal of Machine Learning Research*, 18(154):1–28, 2018.
- S. De Rooij, T. van Erven, P.D. Grünwald, and W.M. Koolen. Follow the leader if you can, hedge if you must. *Journal of Machine Learning Research*, 15(37):1281–1316, 2014.

- J.L. Doob. *Stochastic Processes*. Wiley Publications in Statistics. John Wiley & Sons, 1953.
- D.A Freedman. On tail probabilities for martingales. *The Annals of Probability*, 3(1):100–118, 1975.
- Y. Freund, R.E. Schapire, Y. Singer, and M.K. Warmuth. Using and combining predictors that specialize. In *Proceedings of the 29th annual ACM Symposium on Theory of Computing (STOC’97)*, pages 334–343, 1997.
- A. Garivier, H. Hadiji, P. Ménard, and G. Stoltz. KL-UCB-Switch: Optimal regret bounds for stochastic bandits from both a distribution-dependent and a distribution-free viewpoints, 2019a. Preprint, arXiv:1805.05071.
- A. Garivier, P. Ménard, and G. Stoltz. Explore first, exploit next: The true shape of regret in bandit problems. *Mathematics of Operations Research*, 44(2):377–399, 2019b.
- S. Gerchinovitz and T. Lattimore. Refined lower bounds for adversarial bandits. In *Advances in Neural Information Processing Systems*, pages 1198–1206, 2016.
- J. Honda and A. Takemura. An asymptotically optimal policy for finite support models in the multi-armed bandit problem. *Machine Learning*, 85:361–391, 2011.
- J. Honda and A. Takemura. Non-asymptotic analysis of a new bandit algorithm for semi-bounded rewards. *Journal of Machine Learning Research*, 16(113):3721–3756, 2015.
- J. Kivinen and M.K. Warmuth. Averaging expert predictions. In *Proceedings of the 4th European Conference on Computational Learning Theory (EuroCOLT’99)*, pages 153–167, 1999.
- W.M. Koolen. AdaFTRL. Blog post, Oct. 2016. URL <http://blog.wouterkoolen.info/AdaFTRL/post.html>.
- J. Kwon and V. Perchet. Gains and losses are fundamentally different in regret minimization: The sparse case. *Journal of Machine Learning Research*, 17(227):1–32, 2016.
- T.L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- T. Lattimore. A scale free algorithm for stochastic bandits with bounded kurtosis. In *Advances in Neural Information Processing Systems*, pages 1584–1593, 2017.
- T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- P. Massart. *Concentration Inequalities and Model Selection*, volume XXXIII of *Ecole d’Eté de Probabilités de Saint-Flour*. Springer, 2007. Lectures given in 2003, published in 2007.
- H.B. McMahan. A survey of algorithms and analysis for adaptive online learning. *Journal of Machine Learning Research*, 18(1):3117–3166, 2017.
- Francesco Orabona and Dávid Pál. Scale-free online learning. *Theoretical Computer Science*, 716: 50–69, 2018.
- M.D. Reid, R.M. Frongillo, R.C. Williamson, and N. Mehta. Generalized mixability via entropic duality. In *Proceedings of the 28th Conference on Learning Theory (COLT’15)*, volume 40 of PMLR, pages 1501–1522, 2015.
- Y. Seldin and G. Lugosi. An improved parametrization and analysis of the EXP3++ algorithm for stochastic and adversarial bandits. In *Proceedings of the 30th Annual Conference on Learning Theory (COLT’17)*, volume 65 of PMLR, pages 1743–1759, 2017.



- C.-Y. Wei and H. Luo. More adaptive algorithms for adversarial bandits. In *Proceedings of the 31st Conference On Learning Theory (COLT'18)*, volume 75 of PMLR, pages 1263–1291, 2018.
- J. Zimmert and Y. Seldin. An optimal algorithm for stochastic and adversarial bandits. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTats'20)*, volume 89 of PMLR, pages 467–475, 2019.

# Supplementary material for “Adaptation to the Range in $K$ -Armed Bandits”

by Hédi Hadiji and Gilles Stoltz

## D. Adaptation to the Range for Linear Bandits

To illustrate the generality of the techniques discussed in this paper, we quickly describe how these can be used to obtain range adaptive algorithms for linear bandits. This section is meant for illustration and not for completeness. In particular, we focus on the case of (oblivious) adversarial linear bandits: we refer the reader to Lattimore and Szepesvári [2020, Chapter 27], which we follow closely, for a more thorough description of the setting; we do not describe the application of our techniques to stochastic linear bandits.

**Learning protocol.** A finite action set  $\mathcal{A} \subset \mathbb{R}^d$ , of cardinality  $K$ , is given. (The setting of vanilla  $K$ -armed bandits considered in the rest of the article corresponds to  $\mathcal{A}$  formed by the vertices of the probability simplex of  $\mathbb{R}^K$ .) The environment selects beforehand a sequence  $(y_t)_{t \geq 1}$  of vectors in  $\mathbb{R}^d$  satisfying a boundedness assumption: there exists an interval  $[m, M]$  such that

$$\forall t \geq 1, \quad \forall x \in \mathcal{A}, \quad x^\top y_t \in [m, M]. \quad (22)$$

We assume that the player does not know in advance  $m$  nor  $M$ . To simplify the exposition, we also assume that  $m \leq 0 \leq M$ .

At every time step, the player chooses an action  $X_t \in \mathcal{A}$  and receives and only observes the payoff  $X_t^\top y_t$ . It does not observe  $y_t$  nor the payoffs  $x^\top y_t$  associated with choices  $x \neq X_t$ . The action  $X_t$  is chosen independently at random according to a distribution over  $\mathcal{A}$  denoted by  $p_t = (p_t(a))_{a \in \mathcal{A}}$ .

The expected regret is defined as

$$R_T(y_{1:T}) = \max_{x \in \mathcal{A}} \sum_{t=1}^T x^\top y_t - \mathbb{E} \left[ \sum_{t=1}^T X_t^\top y_t \right].$$

**Estimating the unobserved payoffs.** As in the case of vanilla  $K$ -armed bandits, the key is to estimate unobserved payoffs. We may actually build an estimate  $\hat{y}_t$  of the vectors  $y_t$ , from which we form the estimates  $x^\top \hat{y}_t$ . This estimate takes advantage of the linear structure of the problem.

Fix a distribution  $\pi$  such that the non-negative symmetric matrix

$$M(\pi) \stackrel{\text{def}}{=} \sum_{x \in \mathcal{A}} \pi(x) x x^\top$$

is invertible: such a distribution exists whenever  $\mathcal{A}$  spans  $\mathbb{R}^d$ , which we may assume with no loss of generality; see Lemma 1 below. This distribution  $\pi$  will be used to explore the arms; it is in general not uniform over the arms. For all distributions  $q$  over  $\mathcal{A}$  and all  $\gamma \in (0, 1]$ , the distribution  $p = (1 - \gamma)q + \gamma\pi$  is such that the non-negative symmetric matrix  $M(p)$  is invertible as well (as it is larger than  $\gamma M(\pi)$ , in the sense of the partial inequality  $\succcurlyeq$  over non-negative symmetric matrices). We only use distributions  $p_t$  of this form. We may then define

$$\hat{y}_t = M(p_t)^{-1} X_t X_t^\top y_t \quad (23)$$

---

**Algorithm 2** AdaHedge for adversarial linear bandits
 

---

- 1: **Input:** an exploration distribution  $\pi$  over  $\mathcal{A}$  and exploration rates  $(\gamma_t)_{t \geq 1}$  in  $[0, 1]$
- 2: **Initialization:**  $\eta_1 = +\infty$  and  $q_1$  is the uniform distribution over  $\mathcal{A}$
- 3: **for** rounds  $t = 1, \dots$  **do**
- 4:   Define  $p_t$  by mixing  $q_t$  with  $\pi$  according to
 
$$p_t = (1 - \gamma_t)q_t + \gamma_t\pi$$
- 5:   Draw an arm  $X_t \sim p_t$  (independently at random according to the distribution  $p_t$ )
- 6:   Get and observe the payoff  $X_t^\top y_t$
- 7:   Compute estimates  $x^\top \hat{y}_t$  of all payoffs according to (23)
- 8:   Compute the mixability gap  $\delta_t$  based on the distribution  $q_t$  and on these estimates:

$$\delta_t = \begin{cases} -\sum_{x \in \mathcal{A}} q_t(x) x^\top \hat{y}_t + \frac{1}{\eta_t} \ln \left( \sum_{x \in \mathcal{A}} q_t(x) e^{\eta_t x^\top \hat{y}_t} \right) & \text{if } \eta_t < +\infty \\ -\sum_{x \in \mathcal{A}} q_t(x) x^\top \hat{y}_t + \max_{x \in \mathcal{A}} x^\top \hat{y}_t & \text{if } \eta_t = +\infty \end{cases}$$

- 9:   Compute the learning rate  $\eta_{t+1} = \left( \sum_{s=1}^t \delta_s \right)^{-1} \ln K$
  - 10:   Define  $q_{t+1}$  component-wise as
 
$$q_{t+1}(a) = \exp \left( \eta_{t+1} \sum_{s=1}^t a^\top \hat{y}_s \right) / \sum_{x \in \mathcal{A}} \exp \left( \eta_{t+1} \sum_{s=1}^t x^\top \hat{y}_s \right)$$
  - 11: **end for**
- 

and note that

$$\mathbb{E}[\hat{y}_t | p_t] = M(p_t)^{-1} \underbrace{\left( \sum_{x \in \mathcal{A}} p_t(x) x x^\top y_t \right)}_{=M(p_t)} = y_t; \quad (24)$$

indeed, conditioning on  $p_t$  amounts to integrating over the random choice of  $X_t$  according to  $p_t$ .

**An algorithm adaptive to the unknown range.** When the range is given, a well-known strategy is to use plain exponential weights over actions in  $\mathcal{A}$  with the estimates  $x^\top \hat{y}_t$  to obtain distributions  $q_t$  that are then mixed with  $\pi$  to form the final distributions  $p_t$ . When the range is unknown, we suggest to simply replace plain exponential weights with AdaHedge (the difference lies in the tuning of the rates  $\eta_t$ ), which leads to Algorithm 2. In this algorithm, we refer to rates  $\gamma_t$  as exploration rates (and not as extra-exploration rates as in Algorithm 1) and similarly, to  $\pi$  as the exploration distribution. This is because for adversarial linear bandits, exploration was always required even to get expected results (unlike for  $K$ -armed bandits, see the introduction of Section 4).

The analysis of this algorithm relies on the same ingredients as the ones already encountered in Section 4.3, with the addition of the following lemma, that quantifies the quality of the exploration. This lemma requires that  $\mathcal{A}$  spans  $\mathbb{R}^d$ , which we may assume with no loss of generality (otherwise, we just replace  $\mathbb{R}^d$  by the vector space generated by  $\mathcal{A}$ ).

**Lemma 1** (Lattimore and Szepesvári, 2020, Theorem 21.1). *There exists a distribution  $\pi$  over  $\mathcal{A}$  such that*

$$M(\pi) = \sum_{x \in \mathcal{A}} \pi(x) x x^\top \text{ is invertible} \quad \text{and} \quad \max_{x \in \mathcal{A}} x^\top M(\pi)^{-1} x = d.$$

We are now ready to state the main result of this section. It is the counterpart of Theorem 2; for the sake of simplicity, we only state it for the value  $\alpha = 1/2$ .

**Theorem 5.** *AdaHedge for adversarial linear bandits (Algorithm 2) with the extra-exploration*

$$\gamma_t = \min\left\{1/2, \sqrt{2.5 d(\ln K)t^{-1/2}}\right\}$$

*ensures that for all bounded ranges  $[m, M]$  containing 0, for all oblivious individual sequences  $y_1, y_2, \dots$  satisfying the boundedness condition (22),*

$$R_T(y_{1:T}) \leq 12(M - m)\sqrt{dT \ln K} + 18(M - m)d \ln K.$$

The proof starts by following closely the one of Theorem 2 (provided in Appendix B); the differences are underlined and dealt with in the second part of the proof.

*Proof.* By Reminder 2, since the player plays the AdaHedge strategy over the payoff estimates  $x^\top \hat{y}_t$ , the pre-regret satisfies

$$\max_{x \in \mathcal{A}} \sum_{t=1}^T x^\top \hat{y}_t - \sum_{t=1}^T \sum_{a \in \mathcal{A}} q_t(a) a^\top \hat{y}_t \leq 2\sqrt{V_T \ln K} + M_T \left(2 + \frac{4}{3} \ln K\right)$$

with  $V_T = \sum_{t=1}^T \sum_{x \in \mathcal{A}} q_t(x) (x^\top \hat{y}_t)^2$  and

$$M_T = \max\{x^\top \hat{y}_t : t \leq T \text{ and } x \in \mathcal{A}\} - \min\{x^\top \hat{y}_t : t \leq T \text{ and } x \in \mathcal{A}\}.$$

Since  $\gamma_t \leq 1/2$ , we have  $q_t(x) \leq 2p_t(x)$  for all  $x \in \mathcal{A}$ . We therefore define

$$V'_T = \sum_{t=1}^T \sum_{x \in \mathcal{A}} p_t(x) (x^\top \hat{y}_t)^2$$

and have  $V_t \leq 2V'_T$ . By the tower rule, based on the equality (24), and given that the expectation of a maximum is larger than the maximum of the expectations (for the first inequality), and by the definition of the  $p_t$  (for the second inequality), we have proved so far that

$$\begin{aligned} R_T(y_{1:T}) &\leq \mathbb{E} \left[ \max_{x \in \mathcal{A}} \sum_{t=1}^T x^\top \hat{y}_t - \sum_{t=1}^T \sum_{a \in \mathcal{A}} p_t(a) a^\top \hat{y}_t \right] \\ &\leq \mathbb{E} \left[ \max_{x \in \mathcal{A}} \sum_{t=1}^T x^\top \hat{y}_t - \sum_{t=1}^T \sum_{a \in \mathcal{A}} q_t(a) a^\top \hat{y}_t \right] + \mathbb{E} \left[ \sum_{t=1}^T \gamma_t \sum_{a \in \mathcal{A}} (\pi(a) - q_t(a)) a^\top \hat{y}_t \right] \\ &\leq \mathbb{E} \left[ 2\sqrt{2V'_T \ln K} + M_T \left(2 + \frac{4}{3} \ln K\right) \right] + \underbrace{\sum_{t=1}^T \gamma_t \sum_{a \in \mathcal{A}} (\pi(a) - q_t(a)) a^\top \hat{y}_t}_{\leq (M-m)} \end{aligned}$$

Hence by Jensen's inequality and by the bounds  $\mathbb{E}[V'_T] \leq (M - m)^2 dT$  and  $M_T \leq 2(M - m)d/\gamma_T$  proved below, we finally get

$$\begin{aligned} R_T(y_{1:t}) &\leq 2\sqrt{2\mathbb{E}[V'_T] \ln K} + \mathbb{E}[M_T] \left(2 + \frac{4}{3} \ln K\right) + (M - m) \sum_{t=1}^T \gamma_t \\ &\leq 2\sqrt{2}(M - m)\sqrt{dT \ln K} + \left(2 + \frac{4}{3} \ln K\right) \frac{2(M - m)d}{\gamma_T} + (M - m) \sum_{t=1}^T \gamma_t \\ &\leq 3(M - m)\sqrt{dT \ln K} + 9(M - m) \frac{d \ln K}{\gamma_T} + (M - m) \sum_{t=1}^T \gamma_t. \end{aligned}$$

Replacing the  $\gamma_t$  by their values and using the same bounds as at the beginning of Appendix B yields the claimed result; the factor 12 in the bound comes from

$$3 + \sqrt{10} + 9\sqrt{\frac{2}{5}} \leq 12.$$

We only need to prove the two claimed bounds to complete the proof; they can be extracted from the proof of Theorem 27.1 by Lattimore and Szepesvári [2020] but we provide derivations for the sake of completeness.

*Proof of  $M_T \leq 2(M - m)d/\gamma_T$ .* We fix  $x \in \mathcal{A}$  and  $t \leq T$ . We recall that  $M(p_t)$  and thus  $M(p_t)^{-1}$  are positive definite symmetric matrices. By the Cauchy-Schwarz inequality applied with the norm induced by the positive  $M(p_t)^{-1}$ ,

$$|x^\top M(p_t)^{-1} X_t| \leq \sqrt{x^\top M(p_t)^{-1} x} \sqrt{X_t^\top M(p_t)^{-1} X_t} \leq \max_{x \in \mathcal{A}} \left\{ x^\top M(p_t)^{-1} x \right\}.$$

As indicated right before (23), we have  $M(p_t) \succcurlyeq \gamma_t M(\pi)$  and therefore  $M(p_t)^{-1} \preccurlyeq M(\pi)^{-1}/\gamma_t$ . This entails

$$|x^\top M(p_t)^{-1} X_t| \leq \frac{1}{\gamma_t} \max_{x \in \mathcal{A}} \left\{ x^\top M(\pi)^{-1} x \right\} = \frac{d}{\gamma_t} \leq \frac{d}{\gamma_T},$$

where the equality follows from Lemma 1 and where we used  $\gamma_T \leq \gamma_t$  for the second inequality. Finally, keeping in mind that we assumed  $m \leq 0 \leq M$ ,

$$x^\top \hat{y}_t = \underbrace{x^\top M(p_t)^{-1} X_t}_{\in [-d/\gamma_t, d/\gamma_t]} \underbrace{X_t^\top y_t}_{\in [m, M]} \in \left[ -\frac{d \max\{-m, M\}}{\gamma_T}, \frac{d \max\{-m, M\}}{\gamma_T} \right],$$

from which the bound

$$M_t = 2 \frac{d \max\{-m, M\}}{\gamma_T} \leq \frac{2d(M - m)}{\gamma_T}$$

follows, as desired.

*Proof of  $\mathbb{E}[V'_T] \leq (M - m)^2 dT$ .* Since  $|X_t^\top y_t| \leq \max\{-m, M\} \leq M - m$ , the definition (23) leads to

$$\begin{aligned} (x^\top \hat{y}_t)^2 &= \left( x^\top M(p_t)^{-1} X_t X_t^\top y_t \right)^2 \leq (M - m)^2 \left( x^\top M(p_t)^{-1} X_t \right)^2 \\ &= (M - m)^2 X_t^\top M(p_t)^{-1} x x^\top M(p_t)^{-1} X_t. \end{aligned}$$

Therefore, summing over  $x \in \mathcal{A}$  and using the very definition of  $M(p_t)$ , we get

$$\begin{aligned} \sum_{x \in \mathcal{A}} p_t(x) (x^\top \hat{y}_t)^2 &\leq (M - m)^2 X_t^\top M(p_t)^{-1} \left( \sum_{x \in \mathcal{A}} p_t(x) x x^\top \right) M(p_t)^{-1} X_t \\ &= (M - m)^2 X_t^\top M(p_t)^{-1} X_t = (M - m)^2 \text{Tr} \left( M(p_t)^{-1} X_t X_t^\top \right). \end{aligned}$$

Now, by the linearity of the trace,

$$\mathbb{E} \left[ \text{Tr} \left( M(p_t)^{-1} X_t X_t^\top \right) \right] = \mathbb{E} \left[ \sum_{x \in \mathcal{A}} p_t(x) \text{Tr} \left( M(p_t)^{-1} x x^\top \right) \right] = \mathbb{E} [\text{Tr}(I_d)] = d,$$

where  $I_d$  is the  $d$ -dimensional identity matrix. Collecting all bounds together and summing over  $t$  yields the claimed inequality  $\mathbb{E}[V'_T] \leq (M - m)^2 dT$ .  $\square$

## E. Bandits with Losses (i.e., Known Upper End $M$ on the Range)

In this section, we provide details on the various claims and results hinted at in the paragraph “Known upper end  $M$  on the payoff range” of Section 5.

We will only discuss distribution-free and distribution-dependent upper bounds on the regret, as well as distribution-dependent lower bounds on the regret. This is because the  $(M - m)\sqrt{KT}$  distribution-free regret lower bound of Auer et al. [2002b] holds even in the case when both ends  $m$  and  $M$  of the range are known.

When the upper end  $M$  of the payoff range is known,  $\ln T$  distribution-dependent regret rates are possible and there exists an algorithm achieving the optimal problem-dependent constant (Section E.1). Also,  $\sqrt{KT}$  scale-free distribution-free regret upper bounds may be achieved (Section E.2), which exactly match the distribution-free lower bound. We could not exhibit a strategy that would simultaneously achieve both optimal distribution-dependent and distribution-free regret bounds, unlike what is known in the case of a known payoff range (the KL-UCB-switch strategy by Garivier et al., 2019a) and unlike what we achieved in the main body of the article when adapting to the unknown range or unknown upper end  $M$  on the range. We however conjecture that this should be possible and that, at least, no trade-off exists between the two bounds (i.e., we conjecture that Theorem 1 should not hold).

More precisely, the case considered in this section corresponds to the models  $\mathcal{D}_{-,M}$ , for  $M \in \mathbb{R}$ , defined as

$$\mathcal{D}_{-,M} = \bigcup_{\substack{m \in \mathbb{R}, \\ m < M}} \mathcal{D}_{m,M}.$$

We may adapt Definitions 1 and 2 to define the concepts of distribution-free and distribution-dependent rates for adaptation to the lower end of the range by considering the model  $\mathcal{D}_{-,M}$  therein (just as we did in Remark 1 for the models  $\mathcal{D}_{m,+}$ ).

### E.1. Known $M$ but Unknown $m$ , Part 1: Distribution-Dependent Bounds

The results of this section actually also hold more generally for semi-bounded payoffs, which correspond to the models  $\mathcal{D}_{-\infty,M}$ , for  $M \in \mathbb{R}$ , defined as the sets of probability distributions with a first moment supported on  $(-\infty, M]$ . Note that we have the strict inclusion  $\mathcal{D}_{-,M} \subset \mathcal{D}_{-\infty,M}$  as distributions in  $\mathcal{D}_{-\infty,M}$  are not bounded in general.

The DMED strategy of Honda and Takemura [2015] does achieve a  $\ln T$  distribution-dependent rate for adaptation to the lower end of the range and is even competitive against all bandit problems in  $\mathcal{D}_{-\infty,M}$ . The achieved upper bound is asymptotically optimal as indicated by Reminder 1.

**Reminder 4** (Honda and Takemura, 2015, main theorem). *The regret of the DMED strategy is bounded, for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}_{-\infty,M}$ , by*

$$\limsup_{T \rightarrow \infty} \frac{R_T(\underline{\nu})}{\ln T} \leq \sum_{a=1}^K \frac{\Delta_a}{\mathcal{K}_{\inf}(\nu_a, \mu^*, \mathcal{D}_{-\infty,M})}.$$

The nice and deep result of Reminder 4 implies that from the distribution-dependent point of view, adaptation to the lower end  $m$  of the range is automatic (if such a lower end exists: result holds also when there is no lower bound on the payoffs). Our intuition and understanding for this situation is the following. When the model is  $\mathcal{D}_{m,M}$  for known ends  $m$  and  $M$ , the optimal constant for the  $\ln T$  regret on a bandit problem  $\underline{\nu}$  in  $\mathcal{D}_{m,M}$  is given (see again Reminder 1) by

$$C(\underline{\nu}, m, M) = \sum_{a=1}^K \frac{\Delta_a}{\mathcal{K}_{\inf}(\nu_a, \mu^*, \mathcal{D}_{m,M})}.$$

But it actually turns out, as indicated by Proposition 1 below, that  $C(\underline{\nu}, m, M)$  is independent of  $m$  and equals  $C(\underline{\nu}, -\infty, M)$ .

**Proposition 1.** Fix  $M \in \mathbb{R}$ . For all  $m \leq M$ , for all  $\nu \in \mathcal{D}_{m,M}$  and all  $\mu > E(\nu)$ ,

$$\mathcal{K}_{\inf}(\nu, \mu, \mathcal{D}_{m,M}) = \mathcal{K}_{\inf}(\nu, \mu, \mathcal{D}_{-\infty,M}).$$

*Proof.* The inequality  $\geq$  is immediate, as the right-hand side of the equality is an infimum over the larger set  $\mathcal{D}_{-\infty,M}$ . For the inequality  $\leq$ , we may assume with no loss of generality that  $\mu < M$ , as otherwise, there is no distribution  $\nu'$  neither in  $\mathcal{D}_{m,M}$  nor in  $\mathcal{D}_{-\infty,M}$  with  $E(\nu') > \mu \geq M$ , so that both  $\mathcal{K}_{\inf}$  quantities equal  $+\infty$ .

We fix  $M, m, \nu$  and  $\mu$  as in the statement of the proposition. It suffices to show that in the case  $\mu < M$ , for all  $\nu' \in \mathcal{D}_{-\infty,M}$  with  $E(\nu') > \mu$  and  $\nu \ll \nu'$ , there exists  $\nu'' \in \mathcal{D}_{m,M}$  with  $E(\nu'') > \mu$  and  $\text{KL}(\nu, \nu'') \leq \text{KL}(\nu, \nu')$ . (If  $\nu$  is not absolutely continuous with respect to  $\nu'$ , then  $\text{KL}(\nu, \nu') = +\infty$  and taking  $\nu''$  as the Dirac mass  $\delta_M$  at  $M$  is a suitable choice.) To do so, given such a distribution  $\nu'$ , we first note that  $\nu \ll \nu'$  and  $\nu \in \mathcal{D}_{m,M}$ , i.e.,  $\nu([m, M]) = 1$ , entail that  $\nu'([m, M]) > 0$ , so that we may define the restriction  $\nu'' = \nu'_{[m,M]}$  of  $\nu'$  to  $[m, M]$ ; its density with respect to  $\nu'$  is given by

$$\frac{d\nu''}{d\nu'}(x) = \nu'([m, M])^{-1} \mathbb{1}_{\{x \in [m, M]\}} \quad \nu'\text{-a.s. for all } x \in \mathbb{R}.$$

We have the absolute-continuity chain  $\nu \ll \nu'' \ll \nu'$ , and the Radon-Nykodym derivatives thus defined satisfy

$$\frac{d\nu}{d\nu'}(x) = \frac{d\nu}{d\nu''}(x) \frac{d\nu''}{d\nu'}(x) = \nu'([m, M])^{-1} \frac{d\nu}{d\nu''}(x) \mathbb{1}_{\{x \in [m, M]\}} \quad \nu'\text{-a.s. for all } x \in \mathbb{R}. \quad (25)$$

Moreover  $E(\nu'') \geq E(\nu')$ , and thus  $E(\nu'') > \mu$ , as

$$\begin{aligned} E(\nu') &= \int_{(-\infty, m)} x d\nu'(x) + \int_{[m, M]} x d\nu'(x) \\ &\leq \left(1 - \nu'([m, M])\right)m + \nu'([m, M]) E(\nu'') \leq E(\nu''). \end{aligned}$$

Finally, by (25), which also holds  $\nu$ -almost surely, and the definition of Kullback-Leibler divergences,

$$\begin{aligned} \text{KL}(\nu, \nu') &= \int_{(-\infty, M]} \ln\left(\frac{d\nu}{d\nu'}\right) d\nu = -\ln \nu'([m, M]) + \int_{[m, M]} \ln\left(\frac{d\nu}{d\nu''}\right) d\nu \\ &= -\ln \nu'([m, M]) + \text{KL}(\nu, \nu'') \geq \text{KL}(\nu, \nu''). \end{aligned}$$

This concludes the proof.  $\square$

## E.2. Known $M$ but Unknown $m$ , Part 2: Distribution-Free Bounds

A first observation is that (as in the case of a fully known payoff range) AdaHedge does not require any extra-exploration (i.e., any mixing with the uniform distribution) to achieve a scale-free distribution-free regret bound of order  $(M - m)\sqrt{KT \ln K}$ . This is formally detailed in Appendix E.2.3.

Both this result and the one described next rely on the AdaFTRL methodology of Orabona and Pál [2018], which we recall in Appendix E.2.1. AdaFTRL stands for adaptive follow-the-regularized-leader and it was partially built on and inspired the analysis for AdaHedge, which is a special case of AdaFTRL with entropic regularizer (see De Rooij et al., 2014 for AdaHedge, as well as the earlier analysis by Cesa-Bianchi et al., 2007). Koolen [2016] actually proposes an alternative analysis of AdaFTRL, closer to the AdaHedge formulation, namely, using directly some mixability gaps instead of upper bounds thereon; this is the analysis we recall in Section E.2.1.

The INF strategy of Audibert and Bubeck [2009] can be seen as an instance of FTRL with  $1/2$ -Tsallis entropy, as essentially noted by Audibert et al. [2014]. The INF strategy provides a distribution-free regret bound of order  $\sqrt{KT}$  in case of a known payoff range. Up to some technical issues, which

we could solve, it may be extended to provide a similar scale-free distribution regret bound, which is optimal as it does not contain any superfluous  $\sqrt{\ln K}$  factor. The exact statement to be proved in Appendix E.2.4 is the following: AdaFTRL with  $1/2$ -Tsallis entropy relying on an upper bound  $M$  on the payoffs ensures that for all  $m \in \mathbb{R}$  with  $m \leq M$ , for all oblivious individual sequences  $y_1, y_2, \dots$  in  $[m, M]^K$ , for all  $T \geq 1$ ,

$$R_T(y_{1:T}) \leq 4(M - m)\sqrt{KT} + 2(M - m).$$

We now give a high-level idea of the technical issues that were solved to obtain the bound above. We consider estimates  $\hat{y}_{t,a}$  obtained from (9) by replacing the constant  $C$  therein by the known upper end  $M$ . We however could not simply derive the regret bound from some generic full-information regret guarantee for AdaFTRL with  $1/2$ -Tsallis entropy, as to the best of our knowledge, there are no meaningful full-information regret bounds for Tsallis entropy in the first place, and as these would anyway scale with the effective range of the estimates. We instead provide a more careful analysis exploiting special properties of the estimates:  $\hat{y}_{t,a} = M$  for all  $a \neq A_t$  and  $\hat{y}_{t,A_t} \leq M$ .

We however were unable so far to provide a non-trivial distribution-dependent regret bound for our strategy AdaFTRL with  $1/2$ -Tsallis entropy. Note that there exist  $\mathcal{O}(\ln T)$  bounds for FTRL with  $1/2$ -Tsallis entropy, i.e., with a different tuning of the learning rates (namely,  $\eta_t$  of order  $1/\sqrt{t}$ , but then, the scale-free distribution-free guarantees are lost); see Zimmert and Seldin [2019]. We would have liked to prove such a  $\mathcal{O}(\ln T)$  scale-free distribution-dependent regret bound for AdaFTRL with  $1/2$ -Tsallis entropy (or even achieve a more modest aim like a poly-logarithmic bound), as this seems possible and would have shown with certainty that the trade-off imposed by Theorem 1 does not hold anymore when the upper end  $M$  on the payoff range is known. The techniques of Seldin and Lugosi [2017], which consist in a precise tuning of the extra-exploration in their variant of the Exp3 algorithm of Auer et al. [2002b] together with a gap estimation scheme, or the ones of Zimmert and Seldin [2019] might be helpful to that end. We leave this problem for future research.

### E.2.1. AdaFTRL for Full Information (Reminder of Known Results)

To avoid confusion with the notation used in the main body of the paper, we first describe the considered setting of prediction of oblivious individual sequences with full information.

**Full-information setting.** The game between the player and the environment is actually the same as the one described in Section 4.1, except that the player observes at each step the entire payoff vector, not just the obtained payoff. More formally (and with a different piece of notation  $z$  instead of  $y$ , to better distinguish the two settings), the environment first picks a sequence of payoff vectors  $z_t \in \mathbb{R}^K$ , for all  $t \geq 1$ . Then, in a sequential manner, at every time step  $t$ , the player picks an action  $A_t$ , distributed according to a probability  $p_t$  over the action set  $[K]$ , obtains the payoff  $z_{t,A_t}$ , and observes the entire vector  $z_t$  (i.e., also the payoffs  $z_{t,a}$  corresponding to the actions  $a \neq A_t$ ).

In the sequel, we denote by  $\mathcal{S}$  the simplex of probability distributions over  $[K]$  and we use the short-hand notation, for  $p \in \mathcal{S}$  and  $z \in \mathbb{R}^K$ ,

$$\langle p, z \rangle = \sum_{a \in [K]} p_a z_a.$$

**FTRL (follow-the-regularized-leader).** The FTRL method consists in choosing  $p_t$  according to

$$p_t \in \operatorname{argmin}_{p \in \mathcal{S}: F(p) < +\infty} \left\{ \frac{F(p)}{\eta_t} - \sum_{s=1}^{t-1} \langle p, z_s \rangle \right\},$$

where  $F : \mathbb{R}^K \rightarrow \mathbb{R} \cup \{+\infty\}$  is a convex function, called the regularizer, and  $\eta_t$  is a non-negative learning rate in  $(0, +\infty]$ , which may depend on past observations. The condition  $F(p) < +\infty$  will



always be satisfied for some  $p \in \mathcal{S}$  by the considered regularizers (see below) and is only meant to avoid the undefined  $+\infty / +\infty$  in the case  $\eta_t = +\infty$ . For the sake of concision we will however omit it in the sequel.

Let us give a succinct account of the convex analysis results we use here, following the exposition of Lattimore and Szepesvári [2020, Chapter 26]. Using their terminology, the domain  $\text{Dom } L$  of a convex function  $L : \mathbb{R}^K \rightarrow \mathbb{R} \cup \{+\infty\}$  is the set  $\{x \in \mathbb{R}^K : L(x) < +\infty\}$  of those points where it takes finite values. A convex function  $L : \mathbb{R}^K \rightarrow \mathbb{R} \cup \{+\infty\}$  is said to be Legendre if the interior of its domain  $\text{Int}(\text{Dom } L)$  is non-empty, if  $L$  is strictly convex and differentiable on  $\text{Int}(\text{Dom } L)$ , and if its gradient  $\nabla L$  blows up on the boundary of  $\text{Dom } L$ . The minimizers of Legendre functions may be seen to satisfy the following properties.

**Proposition 2** (Special case of Lattimore and Szepesvári, 2020, Proposition 26.14). *Let  $L$  be a Legendre function and  $A \subseteq \mathbb{R}^d$  be a convex set that intersects  $\text{Int}(\text{Dom } L)$ . Then  $L$  possesses a unique minimizer  $x^*$  over  $A$ , which belongs to  $\text{Int}(\text{Dom } L)$ , therefore ensuring that  $L$  is differentiable at  $x^*$ . Furthermore,*

$$\forall x \in A \cap \text{Dom } L, \quad \langle \nabla L(x^*), x - x^* \rangle \geq 0.$$

Finally, for  $x, y \in \mathbb{R}^d$ , if  $F : \mathbb{R}^K \rightarrow \mathbb{R} \cup \{+\infty\}$  is differentiable at  $y$ , we define the Bregman divergence between  $x$  and  $y$  as

$$B_F(x, y) = F(x) - F(y) - \langle \nabla F(y), x - y \rangle; \quad (26)$$

when  $F$  is convex, we have  $B_F(x, y) \geq 0$  for all  $x \in \mathbb{R}^d$ .

We are now ready to state our first reminder, which is a classical regret bound for FTRL (see, e.g., Lattimore and Szepesvári, 2020, Chapter 28, Exercise 28.12 for references, and McMahan, 2017 for more general versions). It involves the diameter  $D_F$  of the action set (the  $K$ -dimensional simplex  $\mathcal{S}$  in our case):

$$D_F = \max_{p, q \in \mathcal{S}} \{F(p) - F(q)\}.$$

**Reminder 5** (Generic full-information FTRL bound over the simplex). *The FTRL method with a Legendre regularizer  $F$  (of finite diameter  $D_F$ ) and with any rule for picking the learning rates so that they form a non-increasing sequence satisfies the following guarantee: for all sequences  $z_1, z_2, \dots$  of vector payoffs in  $\mathbb{R}^K$ , the regret is bounded by*

$$\begin{aligned} \max_{a \in [K]} \sum_{t=1}^T z_{t,a} - \sum_{t=1}^T \langle p_t, z_t \rangle &\leq \frac{D_F}{\eta_T} + \sum_{t=1}^{T-1} \left( \langle p_t - p_{t+1}, -z_t \rangle - \frac{B_F(p_{t+1}, p_t)}{\eta_t} \right) \\ &\quad + \left( \langle p_T - p^*, -z_T \rangle - \frac{B_F(p^*, p_T)}{\eta_T} \right), \end{aligned} \quad (27)$$

$$\text{where} \quad p^* \in \operatorname{argmax}_{p \in \mathcal{S}} \sum_{t=1}^T \langle p, z_t \rangle$$

and where the regret bound is well defined, thanks to the following observations and conventions: for rounds  $t \geq 1$  where  $\eta_t < +\infty$ , the function  $F$  is indeed differentiable at  $p_t$  so that  $B_F(p_{t+1}, p_t)$  is well defined; for rounds  $t \geq 1$  where  $\eta_t = +\infty$ , we set  $B_F(p_{t+1}, p_t)/\eta_t = 0$  irrespectively of the fact whether  $F$  is differentiable at  $p_t$ .

*Proof of Reminder 5.* Denote by  $S_t$  the cumulative vector payoff up to time  $t \geq 1$ . Fix  $T \geq 1$ . For the sake of concision of the equations, we define  $p_{T+1} = p^*$ , which is a Dirac mass at some arm (that is,  $p_{T+1}$  is not given by FTRL). The regret can therefore be rewritten as

$$\begin{aligned} \max_{a \in [K]} \sum_{t=1}^T z_{t,a} - \sum_{t=1}^T \langle p_t, z_t \rangle &= \max_{p \in \mathcal{S}} \sum_{t=1}^T \langle p, z_t \rangle - \sum_{t=1}^T \langle p_t, z_t \rangle \\ &= \sum_{t=1}^T \langle p_{T+1}, z_t \rangle - \sum_{t=1}^T \langle p_t, z_t \rangle = \sum_{t=1}^T \langle p_t - p_{T+1}, -z_t \rangle. \end{aligned}$$

By summation by parts,

$$\begin{aligned}
 & \sum_{t=1}^T \langle p_t - p_{T+1}, -z_t \rangle \\
 &= \sum_{t=1}^T \sum_{s=t}^T \langle p_s - p_{s+1}, -z_t \rangle = \sum_{s=1}^T \sum_{t=1}^s \langle p_s - p_{s+1}, -z_t \rangle = \sum_{s=1}^T \langle p_s - p_{s+1}, -S_s \rangle \\
 &= \sum_{t=1}^T \langle p_t - p_{t+1}, -z_t \rangle + \sum_{t=1}^T \langle p_t - p_{t+1}, -S_{t-1} \rangle.
 \end{aligned} \tag{28}$$

If  $\eta_t < +\infty$ , then by the optimality condition from Proposition 2 applied to the Legendre function  $L : x \mapsto \eta_t^{-1} F(x) - \langle S_{t-1}, x \rangle$ , we know that  $L$  thus  $F$  are differentiable at  $p_t$  and that

$$\begin{aligned}
 & \langle \eta_t^{-1} \nabla F(p_t) - S_{t-1}, p_{t+1} - p_t \rangle \geq 0, \\
 & \text{that is,} \quad \langle p_t - p_{t+1}, -S_{t-1} \rangle \leq \langle \eta_t^{-1} \nabla F(p_t), p_{t+1} - p_t \rangle.
 \end{aligned}$$

If  $\eta_t = +\infty$ , the previous inequality holds too, as by definition of  $p_t$ , we have  $\langle p_t - p_{t+1}, -S_{t-1} \rangle \leq 0$  and as we set by convention  $\eta_t^{-1} \nabla F(p_t) = 0$  regardless of whether  $F$  is differentiable at  $p_t$  or not. Substituting in (28), we proved so far

$$\sum_{t=1}^T \langle p_t - p_{T+1}, -z_t \rangle \leq \sum_{t=1}^T \langle p_t - p_{t+1}, -z_t \rangle + \langle \eta_t^{-1} \nabla F(p_t), p_{t+1} - p_t \rangle. \tag{29}$$

This inequality can be rewritten in terms of Bregman divergences:

$$\sum_{t=1}^T \langle p_t - p^*, -z_t \rangle \leq \sum_{t=1}^T \left( \langle p_t - p_{t+1}, -z_t \rangle - \frac{B_F(p_{t+1}, p_t)}{\eta_t} \right) + \sum_{t=1}^T \frac{F(p_{t+1}) - F(p_t)}{\eta_t}$$

We now upper bound the second sum in the right-hand side: again by summation by parts, with the convention  $\eta_0 = +\infty$  and  $1/\eta_0 = 0$ :

$$\begin{aligned}
 & \sum_{t=1}^T \frac{F(p_{t+1}) - F(p_t)}{\eta_t} = \sum_{t=1}^T (F(p_{t+1}) - F(p_t)) \sum_{s=1}^t \left( \frac{1}{\eta_s} - \frac{1}{\eta_{s-1}} \right) \\
 &= \sum_{s=1}^T \sum_{t=s}^T (F(p_{t+1}) - F(p_t)) \left( \frac{1}{\eta_s} - \frac{1}{\eta_{s-1}} \right) = \sum_{s=1}^T \underbrace{(F(p_{T+1}) - F(p_s))}_{\leq D_F} \underbrace{\left( \frac{1}{\eta_s} - \frac{1}{\eta_{s-1}} \right)}_{\geq 0} \leq \frac{D_F}{\eta_T},
 \end{aligned}$$

where the final equality is obtained by a telescoping sum, using that the sequence of learning rates is non-increasing.  $\square$

**AdaFTRL, an adaptive version of FTRL.** The AdaFTRL approach consists in tuning the learning rate in a way that scales with the observed data. More precisely, it relies on a quantity called the (generalized) mixability gap, which naturally appears as an upper bound on the summands in the FTRL bound of Remark 5:

$$\delta_t^F \stackrel{\text{def}}{=} \max_{p \in \mathcal{S}} \left\{ \langle p_t - p, -z_t \rangle - \frac{B_F(p, p_t)}{\eta_t} \right\} \geq 0. \tag{30}$$

That mixability gaps are always nonnegative can be seen by taking  $p = p_t$  in the definition. We may further upper bound (27) when it holds by using this mixability gap:

$$\max_{a \in [K]} \sum_{t=1}^T z_{t,a} - \sum_{t=1}^T \langle p_t, z_t \rangle \leq \frac{D_F}{\eta_T} + \sum_{t=1}^T \delta_t^F. \tag{31}$$

The AdaFTRL learning rate balances the two terms in the above regret bound by taking

$$\eta_t = D_F \left/ \sum_{s=1}^{t-1} \delta_s^F \right. \in (0, +\infty] \quad (32)$$

Note that this rule for picking learning rates indeed leads to non-increasing sequences thereof, as the mixability gaps are non-negative. We summarize the discussion above in the theorem stated next, from which subsequent (closed-from) regret bounds will be derived by using the specific properties of the regularizer  $F$  at hand to upper bound the mixability gaps.

**Theorem 6** (AdaFTRL tool box). *Under the assumptions of Reminder 5 and with its conventions, the regret of the FTRL method based on the learning rates (32) satisfies*

$$\max_{a \in [K]} \sum_{t=1}^T z_{t,a} - \sum_{t=1}^T \langle p_t, z_t \rangle \leq 2 \sum_{t=1}^T \delta_t^F \quad (33)$$

where, moreover,

$$\left( \sum_{t=1}^T \delta_t^F \right)^2 = 2D_F \sum_{t=1}^T \frac{\delta_t^F}{\eta_t} + \sum_{t=1}^T (\delta_t^F)^2. \quad (34)$$

*Proof.* Inequality (33) follows from (31) and (32). The equality (34) is obtained by expanding the squared sum,

$$\left( \sum_{t=1}^T \delta_t^F \right)^2 = \sum_{t=1}^T (\delta_t^F)^2 + 2 \sum_{t=1}^T \sum_{s=1}^{t-1} \delta_t^F \delta_s^F = \sum_{t=1}^T (\delta_t^F)^2 + 2 \sum_{t=1}^T \delta_t^F \frac{D_F}{\eta_t}$$

where the final equality is obtained by substituting the definition (32) of  $\eta_t$ .  $\square$

### E.2.2. AdaHedge for Full Information (Reminder of Known Results)

The content of this section is extracted from various sources, out of which the most important is Koolen [2016]. We claim no novelty. This section recalls how the bound for AdaHedge (Reminder 2, for which a direct proof was provided by De Rooij et al. [2014]) can also be seen as a special case of the results of Section E.2.1.

It is well-known (see Freund et al., 1997, Kivinen and Warmuth, 1999, Audibert, 2009), and can be found again by a simple optimization under a linear constraint, that the Hedge weight update corresponds to FTRL with the negentropy as a regularizer:

$$H_{\text{neg}}(p) = \sum_{a=1}^K p_a \ln p_a,$$

with value  $+\infty$  whenever  $p_a = 0$  for some  $a \in [K]$ . That is,

$$\begin{aligned} \operatorname{argmin}_{p \in \mathcal{S}} \left\{ \frac{H_{\text{neg}}(p)}{\eta_t} - \sum_{s=1}^{t-1} \langle p, z_s \rangle \right\} &= \{p_t\} \\ \text{with } p_{t,a} &= \exp \left( \eta_t \sum_{s=1}^{t-1} z_{a,s} \right) \left/ \sum_{k=1}^K \exp \left( \eta_t \sum_{s=1}^{t-1} z_{k,s} \right) \right. \end{aligned} \quad (35)$$

Straightforward calculation show that the regularizer  $H_{\text{neg}}$  is indeed Legendre (see Lattimore and Szepesvári, 2020, Example 26.11) and the  $H_{\text{neg}}$ -diameter of the simplex equals  $D_{H_{\text{neg}}} = \ln K$ . Reminder 5 and Theorem 6 can therefore be applied.

AdaHedge is exactly AdaFTRL with  $H_{\text{neg}}$  as a regularizer. Indeed, the mixability gap (30) can be computed in closed form (as noted by Reid et al., 2015, Lemma 5) and reads in this case:

$$\delta_t^{\text{neg}} = \begin{cases} -\langle p_t, z_t \rangle + \eta_t^{-1} \ln \left( \sum_{a=1}^K p_{t,a} e^{\eta_t z_{t,a}} \right) & \text{if } \eta_t < +\infty, \\ -\langle p_t, z_t \rangle + \max_{a \in [K]} z_{t,a} & \text{if } \eta_t = +\infty. \end{cases} \quad (36)$$

*Proof of the rewriting (36).* When  $\eta_t = +\infty$ , the mixability gap equals, by definition,

$$\delta_t^F = \max_{p \in \mathcal{S}} \{ \langle p_t - p, -z_t \rangle \} = -\langle p_t, z_t \rangle + \max_{p \in \mathcal{S}} \langle p, z_t \rangle = -\langle p_t, z_t \rangle + \max_{a \in [K]} z_{t,a}.$$

For the case  $\eta_t < +\infty$ , the following formula, which is at the heart of the closed-form formula for the Hedge updates (35), will be useful: for any  $S \in \mathbb{R}^d$ ,

$$\min_{p \in \mathcal{S}} \{ H_{\text{neg}}(p) - \langle p, S \rangle \} = \sum_{i=1}^K \frac{e^{S_i}}{\sum_{j=1}^K e^{S_j}} \left( \ln \left( \frac{e^{S_i}}{\sum_{j=1}^K e^{S_j}} \right) - S_i \right) = -\ln \left( \sum_{i=1}^K e^{S_i} \right). \quad (37)$$

When  $\eta_t < +\infty$ , Equation (35) shows that  $p_t$  lies in the interior  $\text{Int}(\mathcal{S})$  of  $\mathcal{S}$ . The Bregman divergence at hand in the definition (30) of the mixability gaps may be simplified into

$$B_F(p, p_t) = H_{\text{neg}}(p) - H_{\text{neg}}(p_t) - \langle \nabla H_{\text{neg}}(p_t), p - p_t \rangle = H_{\text{neg}}(p) - \langle \nabla H_{\text{neg}}(p_t), p \rangle + 1,$$

where the second inequality holds by taking into account the fact that  $H_{\text{neg}}$  is twice differentiable at any  $p \in \text{Int}(\mathcal{S})$ , with

$$\nabla H_{\text{neg}}(p) = (1 + \ln p_i)_{i \in [K]} \quad \text{so that} \quad \langle \nabla H_{\text{neg}}(p), p \rangle = 1 + \sum_{i=1}^K p_i \ln p_i = 1 + H_{\text{neg}}(p).$$

The mixability gaps can therefore be rewritten

$$\begin{aligned} \delta_t^F &= \max_{p \in \mathcal{S}} \left\{ \langle p_t - p, -z_t \rangle - \frac{B_F(p, p_t)}{\eta_t} \right\} \\ &= -\langle p_t, z_t \rangle - \frac{1}{\eta_t} + \frac{1}{\eta_t} \max_{p \in \mathcal{S}} \{ \eta_t \langle p, z_t \rangle - H_{\text{neg}}(p) + \langle \nabla H_{\text{neg}}(p_t), p \rangle \} \\ &= -\langle p_t, z_t \rangle - \frac{1}{\eta_t} - \frac{1}{\eta_t} \min_{p \in \mathcal{S}} \{ H_{\text{neg}}(p) - \langle p, \eta_t z_t + \nabla H_{\text{neg}}(p_t) \rangle \} \end{aligned}$$

Now by (37), specialized with  $S = \eta_t z_t + \nabla H_{\text{neg}}(p_t)$ , we can compute the value of the minimum:

$$\min_{p \in \mathcal{S}} \{ H_{\text{neg}}(p) - \langle p, \eta_t z_t + \nabla H_{\text{neg}}(p_t) \rangle \} = -\ln \left( \sum_{i=1}^K e^{\eta_t z_i + 1 + \ln p_i} \right) = -1 - \ln \left( \sum_{i=1}^K p_i e^{\eta_t z_i} \right).$$

Collecting all equalities together concludes the proof.  $\square$

Reminder 2 is thus a special case of the following bound.

**Theorem 7** (See Lemma 3 and Theorem 6 of De Rooij et al., 2014). *For all sequences of payoffs  $z_{t,a}$  lying in some bounded real-valued interval, denoted by  $[b, B]$ , for all  $T \geq 1$ , the regret of the AdaHedge algorithm with full information, as defined by (35) and (36), satisfies*

$$\begin{aligned} \max_{k \in [K]} \sum_{t=1}^T z_{t,k} - \sum_{t=1}^T \sum_{a=1}^K p_{t,a} z_{t,a} &\leq 2 \sum_{t=1}^T \delta_t^{\text{neg}} \\ \text{where} \quad \sum_{t=1}^T \delta_t^{\text{neg}} &\leq \sqrt{\sum_{t=1}^T \sum_{a=1}^K p_{t,a} \left( z_{t,a} - \sum_{k \in [K]} q_{t,k} z_{t,k} \right)^2} \ln K + (B - b) \left( 1 + \frac{2}{3} \ln K \right), \end{aligned}$$

and AdaHedge does not require the knowledge of  $[b, B]$  to achieve this bound.

The quantities

$$v_t \stackrel{\text{def}}{=} \sum_{a=1}^K p_{t,a} \left( z_{t,a} - \sum_{k \in [K]} q_{t,k} z_{t,k} \right)^2$$

in the bound correspond to the variance of the random variables taking values  $z_{t,a}$  with probability  $p_{t,a}$ ; the variational formula for variances indicates that

$$\sum_{a=1}^K p_{t,a} \left( z_{t,a} - \sum_{k \in [K]} q_{t,k} z_{t,k} \right)^2 = \min_{c \in \mathbb{R}} \sum_{a=1}^K p_{t,a} (z_{t,a} - c)^2,$$

which entails the final bound given as a note in the statement of Reminder 2.

The following formulation of Bernstein's inequality will be useful in the proof of Theorem 7.

**Lemma 2** (Bernstein's inequality tailored to our needs). *Let  $X$  be a random variable in  $[0, 1]$ , with variance denoting by  $\text{Var}(X)$ . Then for all  $\eta > 0$ ,*

$$\frac{\ln(\mathbb{E}[e^{\eta(X - \mathbb{E}[X])}])}{\eta^2} \leq \frac{1}{2} \text{Var}(X) + \frac{1}{3} \frac{\ln(\mathbb{E}[e^{\eta(X - \mathbb{E}[X])}])}{\eta}.$$

*Proof.* Denote by  $\psi_X(\eta) = \ln(\mathbb{E}[e^{\eta(X - \mathbb{E}[X])}])$  the log-moment generating function of  $X$ . A version of Bernstein's inequality with an appropriate control of the moments (as stated by Massart, 2007, Section 2.2.3 and applied to  $X$  with  $c = 1/3$ ) indicates that for all  $\eta \in (0, 3)$ ,

$$\left(1 - \frac{\eta}{3}\right) \psi_X(\eta) \leq \frac{\eta^2}{2} \text{Var}(X).$$

Actually, this inequality also holds for  $\eta \geq 3$  as its left-hand side is non-positive while its right-hand side is nonnegative. The claimed result is derived by rearranging the terms

$$\psi_X(\eta) \leq \frac{\eta^2}{2} \text{Var}(X) + \frac{\eta}{3} \psi_X(\eta)$$

and by dividing both sides by  $\eta^2$ . □

*Proof of Theorem 7.* We apply Theorem 6. To that end, we first bound the mixability gaps. The rewriting (36) (and Jensen's inequality) directly shows that  $0 \leq \delta_t^{\text{neg}} \leq B - b$ . We may also prove the bound

$$\frac{\delta_t^{\text{neg}}}{\eta_t} \leq \frac{v_t}{2} + \frac{1}{3} (B - b) \delta_t^{\text{neg}}. \quad (38)$$

It suffices to do so for  $\eta_t < +\infty$ . Consider the random variable  $X$  taking values  $(z_{t,a} - b)/(B - b)$  with probability  $p_{t,a}$ , for  $a \in \{1, \dots, K\}$ . The mixability gap can be rewritten as

$$\delta_t^{\text{neg}} = \frac{1}{\eta_t} \psi_X(\eta_t (B - b))$$

with the notation of the proof of Lemma 2. The variance of  $X$  equals  $v_t/(B - b)^2$ . Lemma 2 with  $\eta = \eta_t (B - b)$  yields

$$\frac{\delta_t^{\text{neg}}}{\eta_t (B - b)^2} \leq \frac{v_t}{2(B - b)^2} + \frac{\delta_t^{\text{neg}}}{3(B - b)}.$$

from which we obtain (38) by rearranging.

From (34) and (38), we deduce, together with the bound  $(\delta_t^{\text{neg}})^2 \leq (B - b) \delta_t^{\text{neg}}$ , that

$$\left( \sum_{t=1}^T \delta_t^{\text{neg}} \right)^2 \leq (\ln K) \sum_{t=1}^T v_t + (B - b) \left( \frac{2}{3} \ln K + 1 \right) \sum_{t=1}^T \delta_t^{\text{neg}}.$$

Therefore, using the fact that  $x^2 \leq a + bx$  implies  $x \leq \sqrt{a} + b$  for all  $a, b, x \geq 0$ ,

$$\sum_{t=1}^T \delta_t^{\text{neg}} \leq \sqrt{\ln K \sum_{t=1}^T v_t} + (B - b) \left( \frac{2}{3} \ln K + 1 \right),$$

which thanks to (33) concludes the proof of Theorem 7.  $\square$

### E.2.3. AdaHedge with Known Payoff Upper Bound $M$ (Application of Section E.2.2)

We show how to obtain a scale-free distribution-free regret bound of order  $(M - m)\sqrt{KT \ln K}$  with no extra-exploration (including no initial exploration) when an upper bound  $M$  on the payoffs is given to the player. We consider Algorithm 3, where no mixing takes place (unlike in Algorithm 1) and where the probability distributions  $p_t$  are directly computed via an AdaHedge update (no need for intermediate probabilities  $q_t$ ). Note also that we use the estimates (9) with the choice  $C_t = M$ , that is,

$$\hat{y}_{t,a} = \frac{y_{t,a} - M}{p_{t,a}} \mathbb{1}_{\{A_t=a\}} + M. \quad (39)$$

The following observation is key in the analysis below:  $\hat{y}_{t,a} = M$  for all  $a \neq A_t$  and  $\hat{y}_{t,A_t} \leq M$ . We will also use, as in the proof of Theorem 2,

$$\sum_{a=1}^K p_{t,a} \hat{y}_{t,a} = y_{t,A_t}.$$

The performance bound for this simpler algorithm is stated next.

**Theorem 8.** *AdaHedge for  $K$ -armed bandits relying on an upper bound  $M$  on the payoffs (Algorithm 3) ensures that for all  $m \in \mathbb{R}$  with  $m \leq M$ , for all oblivious individual sequences  $y_1, y_2, \dots$  in  $[m, M]^K$ , for all  $T \geq 1$ ,*

$$R_T(y_{1:T}) \leq 2(M - m)\sqrt{KT \ln K} + 2(M - m).$$

The main technical difference with respect to the analysis of Algorithm 1 is that the mixability gaps are directly bounded by the range  $M - m$ . We no longer need to artificially control the size of the estimates (which we did via extra-exploration) to get, in turn, a control of the mixability gaps.

**Lemma 3** (Improved mixability gap bound). *The mixability gaps of AdaHedge for  $K$ -armed bandits relying on an upper bound  $M$  on the payoffs (Algorithm 3) are bounded, for all  $m \in \mathbb{R}$  with  $m \leq M$ , for all oblivious individual sequences  $y_1, y_2, \dots$  in  $[m, M]^K$ , for all  $t \geq 1$ , by*

$$0 \leq \delta_t \leq M - m \quad \text{and} \quad \frac{\delta_t}{\eta_t} \leq \frac{1}{2} p_{t,A_t}^{-1} (M - y_{t,A_t})^2.$$

*Proof.* The fact that  $\delta_t \geq 0$  holds by definition of the gaps and Jensen's inequality. For  $\delta_t \leq M - m$ , the observations after (39) indicate that when  $\eta_t = +\infty$ ,

$$\delta_t = - \sum_{a=1}^K p_{t,a} \hat{y}_{t,a} + \max_{a \in [K]} \hat{y}_{t,a} = M - \hat{y}_{t,A_t},$$

while for  $\eta_t < +\infty$ ,

$$\begin{aligned} \delta_t &= -y_{t,A_t} + \frac{1}{\eta_t} \ln \left( (1 - p_{t,A_t}) e^{\eta_t M} + p_{t,A_t} e^{\eta_t M} e^{\eta_t (y_{t,A_t} - M)/p_{t,A_t}} \right) \\ &\leq M - y_{t,A_t} + \frac{1}{\eta_t} \ln \left( (1 - p_{t,A_t}) + p_{t,A_t} \underbrace{e^{\eta_t (y_{t,A_t} - M)/p_{t,A_t}}}_{\leq 1} \right), \end{aligned}$$

---

**Algorithm 3** AdaHedge for  $K$ -armed bandits, when an upper bound on the payoffs is given
 

---

- 1: **Input:** an upper bound  $M$  on the payoffs
- 2: **AdaHedge initialization:**  $\eta_1 = +\infty$  and  $p_1 = (1/K, \dots, 1/K)$
- 3: **for** rounds  $t = 1, 2, \dots$  **do**
- 4:   Draw an arm  $A_t \sim p_t$  (independently at random according to the distribution  $p_t$ )
- 5:   Get and observe the payoff  $y_{t,A_t}$
- 6:   Compute the estimates of all payoffs

$$\hat{y}_{t,a} = \frac{y_{t,a} - M}{p_{t,a}} \mathbb{1}_{\{A_t=a\}} + M$$

- 7:   Compute the mixability gap  $\delta_t$  based on the distribution  $p_t$  and on these estimates:

$$\delta_t = \begin{cases} -\sum_{a=1}^K p_{t,a} \hat{y}_{t,a} + \frac{1}{\eta_t} \ln \left( \sum_{a=1}^K p_{t,a} e^{\eta_t \hat{y}_{t,a}} \right) & \text{if } \eta_t < +\infty \\ -\sum_{a=1}^K p_{t,a} \hat{y}_{t,a} + \max_{a \in [K]} \hat{y}_{t,a} & \text{if } \eta_t = +\infty \end{cases}$$

- 8:   Compute the learning rate  $\eta_{t+1} = \left( \sum_{s=1}^t \delta_s \right)^{-1} \ln K$
- 9:   Define  $p_{t+1}$  component-wise as

$$p_{t+1,a} = \exp \left( \eta_{t+1} \sum_{s=1}^t \hat{y}_{a,s} \right) / \sum_{k=1}^K \exp \left( \eta_{t+1} \sum_{s=1}^t \hat{y}_{k,s} \right)$$

- 10: **end for**

---

which entails  $\delta_t \leq M - y_{t,A_t} \leq M - m$ .

Furthermore, in the case  $\eta_t < +\infty$ , using the inequality  $e^{-x} \leq 1 - x + x^2/2$  valid for  $x \geq 0$ , followed by the inequality  $\ln(1 + u) \leq u$ , valid for all  $u > -1$ , we get

$$\delta_t \leq M - \hat{y}_{t,A_t} + \frac{1}{\eta_t} \ln \left( \underbrace{1 - p_{A_t,t} + p_{A_t,t}}_{=1} \underbrace{-\eta_t(M - y_{t,A_t}) + \eta_t^2 \frac{(M - y_{t,A_t})^2}{2p_{A_t,t}}}_{=u} \right) \leq \eta_t \frac{(M - y_{t,A_t})^2}{2p_{t,A_t}}.$$

The second inequality is trivial in case  $\eta_t = +\infty$ , as  $\delta_t/\eta_t = 0$ . □

We are now ready to prove Theorem 8.

*Proof of Theorem 8.* As indicated in Section E.2.2, AdaHedge is a special case of AdaFTRL and the bound of Theorem 6 is applicable.

Equation (34) and Lemma 3, which entails in particular that  $\delta_t^2 \leq (M - m)\delta_t$ , yield

$$\left( \sum_{t=1}^T \delta_t \right)^2 = 2(\ln K) \sum_{t=1}^T \frac{\delta_t}{\eta_t} + \sum_{t=1}^T (\delta_t)^2 \leq (\ln K) \sum_{t=1}^T p_{t,A_t}^{-1} (M - y_{t,A_t})^2 + (M - m) \sum_{t=1}^T \delta_t,$$

which, through the fact that  $x^2 \leq a + bx$  implies  $x \leq \sqrt{a} + b$  for all  $a, b, x \geq 0$ , leads in turn to

$$\sum_{t=1}^T \delta_t \leq \sqrt{\sum_{t=1}^T p_{t,A_t}^{-1} (M - \hat{y}_{t,A_t})^2 \ln K} + (M - m).$$

Therefore, Equation (33) guarantees that

$$\max_{k \in [K]} \sum_{t=1}^T \hat{y}_{t,k} - \underbrace{\sum_{t=1}^T \sum_{a=1}^K p_{t,a} \hat{y}_{t,a}}_{=y_{t,A_t}} \leq 2 \sqrt{\sum_{t=1}^T p_{t,A_t}^{-1} (M - \hat{y}_{t,A_t})^2 \ln K} + 2(M - m). \quad (40)$$

We conclude the proof by integrating the inequality above and using Jensen's inequality, exactly as in the proof of Theorem 2. Indeed, Equation (15) therein indicates that

$$R_T(y_{1:T}) = \max_{k \in [K]} \sum_{t=1}^T y_{t,k} - \mathbb{E} \left[ \sum_{t=1}^T y_{t,A_t} \right] \leq \mathbb{E} \left[ \max_{k \in [K]} \sum_{t=1}^T \hat{y}_{t,k} - \sum_{t=1}^T y_{t,A_t} \right]$$

and, by the same manipulations as in (17) and in the equation that follows it,

$$\begin{aligned} \mathbb{E} \left[ \sqrt{\sum_{t=1}^T p_{t,A_t}^{-1} (M - \hat{y}_{t,A_t})^2 \ln K} \right] &\leq \sqrt{\mathbb{E} \left[ \sum_{t=1}^T p_{t,A_t}^{-1} (M - y_{t,A_t})^2 \ln K \right]} \\ &= \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \sum_{a=1}^K (M - y_{t,a})^2 \ln K \right]} \leq (M - m) \sqrt{KT \ln K} \end{aligned}$$

The claimed result is obtained by collecting all bounds together.  $\square$

#### E.2.4. AdaFTRL with Tsallis Entropy in the Case of a Known Payoff Upper Bound $M$

In this section we describe how the AdaHedge learning rate scheme can be used in the FTRL framework with a different regularizer, namely Tsallis entropy, to improve the scale-free distribution-free regret bound into a bound of optimal order  $(M - m)\sqrt{KT}$ , i.e., without any superfluous  $\sqrt{\ln K}$  factor.

**Tsallis entropy.** We focus on the (rescaled)  $1/2$ -Tsallis entropy, which is defined by

$$H_{1/2}(p) = - \sum_{a=1}^K 2\sqrt{p_a}.$$

This regularizer is Legendre over the domain  $[0, +\infty)^K$  (see Lattimore and Szepesvári, 2020, Example 26.10). Its diameter equals

$$D_{H_{1/2}} = \max_{p \in \mathcal{S}} H_{1/2}(p) - \min_{q \in \mathcal{S}} H_{1/2}(q) = -2 - (-2\sqrt{K}) = 2(\sqrt{K} - 1), \quad (41)$$

as for all  $p \in \mathcal{S}$ , we have (by concavity of the square root for the right-most inequality)

$$1 \leq \sum_{a=1}^K p_a \leq \sum_{a=1}^K \sqrt{p_a} \leq \sqrt{K},$$

where 1 is achieved with  $p = (1, 0, \dots, 0)$  and  $\sqrt{K}$  with the uniform distribution.

The function  $H_{1/2}$  is differentiable at all  $q \in (0, +\infty)^K$ , with  $\nabla H_{1/2}(q) = (-1/\sqrt{q_a})_{a \in [K]}$ . The Bregman divergence associated with  $H_{1/2}$  equals, for  $p, q \in \mathcal{S}$  such that  $q_a > 0$  for all  $a$ :

$$\begin{aligned} B_{H_{1/2}}(p, q) &= -2 \sum_{a=1}^K \sqrt{p_a} + 2 \sum_{a=1}^K \sqrt{q_a} + \sum_{a=1}^K \frac{1}{\sqrt{q_a}} (p_a - q_a) \\ &= -2 \sum_{a=1}^K \frac{\sqrt{p_a} - \sqrt{q_a}}{2\sqrt{q_a}} (2\sqrt{q_a} - (\sqrt{p_a} + \sqrt{q_a})) = \sum_{a=1}^K \frac{(\sqrt{p_a} - \sqrt{q_a})^2}{\sqrt{q_a}}. \end{aligned}$$



**AdaFTRL with 1/2-Tsallis entropy.** We consider FTRL with the 1/2-Tsallis entropy on the estimated losses (39):

$$p_t \in \operatorname{argmin}_{p \in \mathcal{S}} \left\{ \frac{H_{1/2}(p)}{\eta_t} - \sum_{s=1}^{t-1} \langle p, \hat{y}_s \rangle \right\} = \operatorname{argmin}_{p \in \mathcal{S}} \left\{ -\frac{1}{\eta_t} \sum_{a=1}^K 2\sqrt{p_a} - \sum_{a=1}^K p_a \sum_{s=1}^{t-1} \hat{y}_{s,a} \right\}.$$

FTRL with the 1/2-Tsallis entropy was essentially introduced by Audibert and Bubeck [2009] to get rid of a  $\sqrt{\ln K}$  factor in the distribution-free regret bound of  $K$ -armed adversarial bandits (with known payoff range). It was later noted by Audibert et al. [2014] that it actually is an instance of mirror descent with Tsallis entropy as a regularizer. More recently, Zimmert and Seldin [2019] showed that this regularizer can obtain quasi-optimal regret bounds for both stochastic and adversarial rewards.

We more precisely consider AdaFTRL with the 1/2-Tsallis, that is, we compute the learning rates  $\eta_t$  based on the mixability gaps (30); see Algorithm 4. We denote by  $\delta_t^{\text{Ts}}$  the mixability gaps (30).

**On the implementation.** For Tsallis entropy, the optimization problems involved in the computation of the updates  $p_t$  and of the mixability gaps  $\delta_t^{\text{Ts}}$  admit a (semi-)explicit formula. Indeed,  $p_t$  can be computed thanks to the formula, for all  $z \in \mathbb{R}^K$ ,

$$\operatorname{argmin}_{p \in \mathcal{S}} \{H_{1/2}(p) - \langle p, z \rangle\} = \operatorname{argmax}_{p \in \mathcal{S}} \left\{ \langle p, z \rangle + \sum_{a=1}^K 2\sqrt{p_a} \right\} = \left( \frac{1}{(c(z) - z_a)^2} \right)_{a \in K}, \quad (42)$$

where  $c(z)$  is an implicit normalization constant, such that the vector lies in the simplex  $\mathcal{S}$  and  $c(z) > z_a$  for all  $a \in [K]$ . This constant  $c(z)$  is in fact the Lagrange multiplier associated with the constraint  $p_1 + \dots + p_K = 1$ . See Zimmert and Seldin [2019] for more details on how to compute  $c(z)$  efficiently, see also Audibert et al. [2014]. To compute the mixability gap, rewrite

$$\begin{aligned} \delta_t^{\text{Ts}} &= \max_{p \in \mathcal{S}} \left\{ \langle p_t - p, -\hat{y}_t \rangle - \frac{H_{1/2}(p) - H_{1/2}(p_t) - \langle \nabla H_{1/2}(p_t), p - p_t \rangle}{\eta_t} \right\} \\ &= \langle p_t, -\hat{y}_t \rangle + \frac{H_{1/2}(p_t)}{\eta_t} - \frac{\langle \nabla H_{1/2}(p_t), p_t \rangle}{\eta_t} + \frac{1}{\eta_t} \max_{p \in \mathcal{S}} \left\{ \langle p, \nabla H_{1/2}(p_t) + \eta_t \hat{y}_t \rangle - H_{1/2}(p) \right\}, \end{aligned} \quad (43)$$

where the maximum in the left-most side of these equalities can be computed efficiently, thanks to (42).

**Analysis of the algorithm.** We provide the following performance bound.

**Theorem 9.** *AdaFTRL with 1/2-Tsallis entropy for  $K$ -armed bandits relying on an upper bound  $M$  on the payoffs (Algorithm 4) ensures that for all  $m \in \mathbb{R}$  with  $m \leq M$ , for all oblivious individual sequences  $y_1, y_2, \dots$  in  $[m, M]^K$ , for all  $T \geq 1$ ,*

$$R_T(y_{1:T}) \leq 4(M - m)\sqrt{KT} + 2(M - m).$$

As in Section E.2.3, the proof scheme is a combination of the AdaFTRL bound of Theorem 6 (which is indeed applicable), together with an improved bound on the mixability gap that exploits the specific shape of the estimates. This bound is stated in the next lemma, which is much similar to Lemma 3.

**Lemma 4.** *The mixability gaps of AdaFTRL with Tsallis entropy for  $K$ -armed bandits relying on an upper bound  $M$  on the payoffs (Algorithm 4) are bounded, for all  $m \in \mathbb{R}$  with  $m \leq M$ , for all oblivious individual sequences  $y_1, y_2, \dots$  in  $[m, M]^K$ , for all  $t \geq 1$ , by*

$$0 \leq \delta_t^{\text{Ts}} \leq M - m \quad \text{and} \quad \frac{\delta_t^{\text{Ts}}}{\eta_t} \leq p_{t,A_t}^{-1/2} (M - y_{t,A_t})^2.$$

The proof of Lemma 4 is postponed to the end of this section and we now proceed with the proof of Theorem 9.

---

**Algorithm 4** AdaFTRL with Tsallis entropy for  $K$ -armed bandits, when an upper bound on the payoffs is given

---

- 1: **Input:** an upper bound  $M$  on the payoffs
- 2: **Initialization:**  $\eta_1 = +\infty$  and  $p_1 = (1/K, \dots, 1/K)$
- 3: **for** rounds  $t = 1, 2, \dots$  **do**
- 4:   Draw an arm  $A_t \sim p_t$  (independently at random according to the distribution  $p_t$ )
- 5:   Get and observe the payoff  $y_{t,A_t}$
- 6:   Compute the estimates of all payoffs

$$\hat{y}_{t,a} = \frac{y_{t,a} - M}{p_{t,a}} \mathbb{1}_{\{A_t=a\}} + M$$

- 7:   Compute the mixability gap  $\delta_t^{\text{Ts}}$  based on the distribution  $p_t$  and on these estimates, e.g., using the efficient implementation stated around (43):

$$\delta_t^{\text{Ts}} = \max_{p \in \mathcal{S}} \left\{ \langle p_t - p, -\hat{y}_t \rangle - \frac{B_{H_{1/2}}(p, p_t)}{\eta_t} \right\}$$

- 8:   Compute the learning rate  $\eta_{t+1} = 2 \left( \sum_{s=1}^t \delta_s^{\text{Ts}} \right)^{-1} (\sqrt{K} - 1)$
- 9:   Define  $p_{t+1}$  as

$$p_{t+1} \in \operatorname{argmin}_{p \in \mathcal{S}} \left\{ - \sum_{a=1}^K p_a \sum_{s=1}^t \hat{y}_{s,a} - \frac{1}{\eta_{t+1}} \sum_{a=1}^K 2\sqrt{p_a} \right\},$$

where an efficient implementation is provided by, e.g., (42)

10: **end for**

---

*Proof of Theorem 9.* The structure of the proof is much similar to the one of Theorem 8, which is why we only sketch our arguments. The bound of Theorem 6 is applicable. We use Lemma 4 with (34) to see that

$$\left( \sum_{t=1}^T \delta_t^{\text{Ts}} \right)^2 \leq 2D_{H_{1/2}} \sum_{t=1}^T p_{t,A_t}^{-1/2} (M - y_{t,A_t})^2 + (M - m) \sum_{t=1}^T \delta_t^{\text{Ts}}. \quad (44)$$

Again, using the fact that for all  $a, b, x \geq 0$ , the inequality  $x^2 \leq a + bx$  implies  $x \leq \sqrt{a} + b$ :

$$\sum_{t=1}^T \delta_t^{\text{Ts}} \leq \sqrt{2D_{H_{1/2}} \sum_{t=1}^T p_{t,A_t}^{-1/2} (M - y_{t,A_t})^2 + (M - m)} \quad (45)$$

By (33), by taking expectations, and by Jensen's inequality:

$$R_T(y_{1:T}) \leq 2\mathbb{E} \left[ \sum_{t=1}^T \delta_t^{\text{Ts}} \right] \leq 2\sqrt{2D_{H_{1/2}} \sum_{t=1}^T \mathbb{E} \left[ p_{t,A_t}^{-1/2} (M - y_{t,A_t})^2 \right]} + 2(M - m). \quad (46)$$

We conclude by observing that for all  $t$ , by definition of the payoff estimates,

$$\begin{aligned} \mathbb{E} \left[ p_{t,A_t}^{-1/2} (M - y_{t,A_t})^2 \right] &= \mathbb{E} \left[ \sum_{a=1}^K p_{t,a} p_{t,a}^{-1/2} (M - y_{t,a})^2 \right] \leq (M - m)^2 \mathbb{E} \left[ \sum_{a=1}^K \sqrt{p_{t,a}} \right] \\ &\leq (M - m)^2 \sqrt{K}, \end{aligned}$$

where the last inequality follows from the concavity of the square root. The final claim is obtained by bounding the diameter  $D_{H_{1/2}}$  by  $2\sqrt{K}$ .  $\square$

We conclude this section by providing a proof of Lemma 4.

*Proof of Lemma 4.* The fact that  $\delta_t^{\text{Ts}} \geq 0$  holds actually for all regularizers and can be seen from the definition (30) with  $p = p_t$ . For the inequality  $\delta_t^{\text{Ts}} \leq M - m$ , we start with elementary manipulations of the definition of the mixability gap (30). Denoting by  $\vec{M}$  the vector with coordinates  $(M, \dots, M)$  and noting that  $\langle p_t - q, \vec{M} \rangle = 0$  for all  $q \in \mathcal{S}$ , we have

$$\delta_t^{\text{Ts}} = \max_{q \in \mathcal{S}} \left\{ \langle p_t - q, -\hat{y}_t \rangle - \frac{B_{H_{1/2}}(q, p_t)}{\eta_t} \right\} = \max_{q \in \mathcal{S}} \left\{ \langle p_t - q, \vec{M} - \hat{y}_t \rangle - \frac{B_{H_{1/2}}(q, p_t)}{\eta_t} \right\}. \quad (47)$$

Since all the coordinates of  $\vec{M} - \hat{y}_t$  are non-negative and by non-negativity of the Bregman divergence, this implies that

$$\delta_t^{\text{Ts}} \leq \langle p_t, \vec{M} - \hat{y}_t \rangle = M - y_{A_t, t} \leq M - m.$$

We now prove the second inequality; we may assume that  $\eta_t < +\infty$ , as the bound holds trivially otherwise. By Proposition 2 (and by calculations similar to the ones performed in the proof of Reminder 5) the maximum in the rewriting (47) of  $\delta_t^{\text{Ts}}$  is achieved on the interior of the domain of  $H_{1/2}$ , which equals  $(0, +\infty)^K$ , thus in the interior of  $\mathcal{S}$ . We therefore only need to prove that

$$\forall q \in \text{Int}(\mathcal{S}), \quad \langle p_t - q, \vec{M} - \hat{y}_t \rangle - \frac{B_{H_{1/2}}(q, p_t)}{\eta_t} \leq \eta_t p_{t, A_t}^{-1/2} (M - y_{t, A_t})^2. \quad (48)$$

We fix such a  $q \in \text{Int}(\mathcal{S})$ , i.e., such that  $q_a > 0$  for all  $a$ . We consider two cases. First, if  $q_{A_t} \geq p_{t, A_t}$ , then, given the observations made after (39),

$$\langle p_t - q, \vec{M} - \hat{y}_t \rangle - \frac{B_{H_{1/2}}(q, p_t)}{\eta_t} = \underbrace{\left( \frac{M - y_{t, A_t}}{p_{t, A_t}} \right)}_{\geq 0} \underbrace{(p_{t, A_t} - q_{A_t})}_{\leq 0} - \frac{B_{H_{1/2}}(q, p_t)}{\eta_t} \leq 0.$$

Otherwise, when  $q_{A_t} < p_{t, A_t}$ , a standard way of bounding the mixability gap, detailed below, indicates that

$$\langle p_t - q, \vec{M} - \hat{y}_t \rangle - \frac{B_{H_{1/2}}(q, p_t)}{\eta_t} \leq \frac{\eta_t}{2} \left\langle \vec{M} - \hat{y}_t, \nabla^2 H_{1/2}(z)^{-1} (\vec{M} - \hat{y}_t) \right\rangle, \quad (49)$$

where  $z$  is some probability distribution of the open segment  $\text{Seg}(q, p_t)$  between  $q$  and  $p_t$ , and where  $\nabla^2 H_{1/2}(z)^{-1}$  denotes the inverse of the positive definite Hessian of  $H_{1/2}$  at  $z$ . Since at  $w \in (0, +\infty)^K$ , the function  $H_{1/2}$  is indeed twice differentiable, with

$$\nabla H_{1/2}(w) = (-w_a^{-1/2})_{a \in [K]} \quad \text{and} \quad \nabla^2 H_{1/2}(w) = \text{Diag}(w_a^{-3/2}/2)_{a \in [K]},$$

we have  $\nabla^2 H_{1/2}(z)^{-1} = \text{Diag}(2z_a^{3/2})_{a \in [K]}$ . We substitute this value into (49) and recall that the vector  $\vec{M} - \hat{y}_t$  has null coordinates except for its  $A_t$ -th coordinate:

$$\frac{\eta_t}{2} \left\langle \vec{M} - \hat{y}_t, \nabla^2 H_{1/2}(z)^{-1} (\vec{M} - \hat{y}_t) \right\rangle = \eta_t z_{A_t}^{3/2} (M - \hat{y}_{t, A_t})^2.$$

Finally, remember that  $z$  lies in the open segment  $\text{Seg}(q, p_t)$  and that we assumed  $q_{A_t} < p_{t, A_t}$ ; we thus also have  $z_{A_t} < p_{t, A_t}$ . As a consequence, using the very definition of  $\hat{y}_{t, A_t}$ ,

$$\eta_t z_{A_t}^{3/2} (M - \hat{y}_{t, A_t})^2 \leq \eta_t p_{t, A_t}^{3/2} (M - \hat{y}_{t, A_t})^2 = \eta_t p_{t, A_t}^{-1/2} (M - y_{t, A_t})^2.$$

Therefore, in all cases, that is, whether  $q_{A_t} \geq p_{t,A_t}$  or  $q_{A_t} < p_{t,A_t}$ , the bound (48) is obtained. It only remains to prove the standard inequality (49).

This inequality is essentially stated as Theorem 26.13 in Lattimore and Szepesvári [2020] but we provide a proof for the sake of completeness. As we assumed that  $\eta_t < +\infty$ , we have (as above, by Proposition 2) that  $p_t$  lies in the interior of  $\mathcal{S}$ . In particular, as both  $p_t$  and  $q$  are in the interior of  $\mathcal{S}$ , the function  $H_{1/2}$  is  $\mathcal{C}^2$  over the closed segment  $\overline{\text{Seg}}(q, p_t)$  between  $q$  and  $p_t$ . Therefore, by the mean-value theorem, there exists  $z$  in the open segment  $\text{Seg}(q, p_t)$  such that

$$\underbrace{H_{1/2}(q) - H_{1/2}(p_t) - \langle \nabla H_{1/2}(p_t), q - p_t \rangle}_{=B_{H_{1/2}}(q, p_t)} = \frac{1}{2} \langle q - p_t, \nabla^2 H_{1/2}(z) (q - p_t) \rangle.$$

It is useful to introduce the standard notation from convex analysis for the local norm (which is indeed a norm because the Hessian is positive definite):

$$\|q - p_t\|_{\nabla^2 H_{1/2}(z)}^2 \stackrel{\text{def}}{=} \langle q - p_t, \nabla^2 H_{1/2}(z) (q - p_t) \rangle.$$

We therefore have so far the rewriting:

$$-\frac{B_{H_{1/2}}(q, p_t)}{\eta_t} = -\frac{1}{2\eta_t} \langle q - p_t, \nabla^2 H_{1/2}(z) (q - p_t) \rangle.$$

Now, by the Cauchy-Schwarz inequality,

$$\begin{aligned} \langle p_t - q, \vec{M} - \hat{y}_t \rangle &= \langle \nabla^2 H_{1/2}(z)^{1/2} (p_t - q), \nabla^2 H_{1/2}(z)^{-1/2} (\vec{M} - \hat{y}_t) \rangle \\ &\leq \|p_t - q\|_{\nabla^2 H_{1/2}(z)} \|\vec{M} - \hat{y}_t\|_{\nabla^2 H_{1/2}(z)^{-1}}. \end{aligned}$$

Combining the rewriting and the bound above, we get

$$\begin{aligned} \langle p_t - q, M - \hat{y}_t \rangle - \frac{B_{H_{1/2}}(q, p_t)}{\eta_t} &\leq \|p_t - q\|_{\nabla^2 H_{1/2}(z)} \|\vec{M} - \hat{y}_t\|_{\nabla^2 H_{1/2}(z)^{-1}} - \frac{1}{2\eta_t} \|q - p_t\|_{\nabla^2 H_{1/2}(z)}^2 \\ &\leq \frac{\eta_t}{2} \|\vec{M} - \hat{y}_t\|_{\nabla^2 H_{1/2}(z)^{-1}}^2, \end{aligned}$$

where we used  $ab - b^2/2 \leq a^2/2$  to get the second inequality. This is exactly (49).  $\square$